# Artificial Intelligence Basics

Apress®

# ARTIFICIAL INTELLIGENCE BASICS

**A NON-TECHNICAL INTRODUCTION**

*Tom Taulli*

Apress®

*Artificial Intelligence Basics: A Non-Technical Introduction*

Tom Taulli
Monrovia, CA, USA

# Contents

# About the Author

**Tom Taulli** has been developing software since the 1980s. In college, he started his first company, which focused on the development of e-learning systems. He created other companies as well, including Hypermart.net that was sold to InfoSpace in 1996. Along the way, Tom has written columns for online publications such as businessweek.com, techweb.com, and Bloomberg.com. He also writes posts on artificial intelligence for Forbes.com and is the advisor to various companies in the AI space. You can reach Tom on Twitter (@ttaulli) or through his web site (www.taulli.com).

# Foreword

As this book demonstrates, the adoption of artificial intelligence (AI) will be a major inflection point in human history. Like other similarly groundbreaking technologies, how it's administered and who has access to it will shape society for generations to come. However, AI stands out from the other transformative technologies of the nineteenth and twentieth centuries—think the steam engine, the electrical grid, genomics, computers, and the Internet—because it doesn't depend exclusively on critically expensive physical infrastructure to enable adoption; after all, many of its benefits can be delivered through existing hardware we all carry around in our pockets. Instead, the fundamental limiting factor when it comes to the mass adoption of AI technology is our shared intellectual infrastructure: education, understanding, and vision.

This is a crucial difference because, if handled correctly, AI can act as a sweeping democratizing force. It has and will eliminate from our lives the drudgery of the past and free up a tremendous amount of human energy and capital. But that "if" is far from certain. AI executed irresponsibly has the power to destabilize large parts of the world economy by causing, as many people fear, a shrinking workforce, reduced purchasing power for the middle class, and an economy without a wide and stable base fueled by an endless debt spiral.

However, before we succumb to pessimism on AI, we should take a look back. Historic though AI's transformative capacity may be—and it is historic—these same issues are and have been at play in the economic landscape for decades, even centuries. AI is, after all, an extension of a trend toward automation that has been at play since Henry Ford. In fact, Zoho itself was born from the tension between automation and egalitarian economic principles. Back in the early 2000s, we came to a realization that has shaped our approach to technology: regular people—small business owners, here and abroad—should have access to the same advanced business automations that the Fortune 500 companies have; otherwise, a huge swath of the population will be locked out of the economy.

At the time, powerful digital software was almost unanimously gated behind rigid contracts, exorbitant fee structures, and complicated on-premise implementations. Big companies could shoulder the burden of such systems, while smaller operators were locked out, putting them at a tremendous disadvantage. We sought to disrupt that by opening up the promise of technology to wider and wider audiences. Over the last two decades, we've endeavored to

increase the value of our products without increasing the price by tapping into the scalability of cloud technology. Our goal is to empower people at all levels of society by pushing down the price of business software while expanding the power of the tools. Access to capital shouldn't limit success; businesses should rise or fall based on the strength of their vision for the future.

Viewed this way, AI is the fulfillment of the promise of technology. It frees people from the constraints of time by enabling them to offload tedious or unpleasant rote labor. It helps them identify patterns at microscopic and macroscopic scales, which humans are not naturally well suited to perceive. It can forecast problems, and it can correct errors. It can save money, time, and even lives.

Seeking to democratize these benefits just as we did for general business software, Zoho has threaded AI throughout our suite of apps. We spent the last six years quietly developing our own internal AI technology, built on the bedrock of our own principles. The result is Zia, an AI assistant who is smart, but not clever. This is a crucial distinction. A smart system has the information and functionality to empower the unique vision and intuition of an active operator. A clever system obfuscates the internal workings of the process, reducing the human to a passive user who simply consumes the insights provided by the machine. AI should be a tool to be wielded, not a lens through which we view the world. To steer such a powerful tool, we must be equipped with the knowledge to understand and operate it without eroding the human quality of our human systems.

The need to stay current on this technology is exactly why a book like *Artificial Intelligence Basics* is so important in today's world. It is the intellectual infrastructure that will enable people—regular people—to tap into the power of AI. Without these kinds of initiatives, AI will tip the balance of power in favor of big companies with big budgets. It's crucial that the general population equip themselves with the skills to understand AI systems, because these systems will increasingly define how we interact with and navigate through the world. Soon, the information contained in this book won't be merely a topic of interest; it will be a prerequisite for participation in the modern economy.

This is how the average person can enjoy the fruits of the AI revolution. In the years to come, how we define work and which activities carry economic value will change. We have to embrace the fact that the future of work may be as foreign to us as a desk job would be to our distant ancestors. But we have to—and should—have faith in the human capacity to innovate new forms of work, even if that work doesn't look like the work we're familiar with. But the first step, before everything else, is to learn more about this new, exciting, and fundamentally democratizing technology.

—Sridhar Vembu, co-founder and CEO of Zoho

# Introduction

On the face of it, the Uber app is simple. With just a couple clicks, you can hail a driver within a few minutes.

But behind the scenes, there is an advanced technology platform, which relies heavily on artificial intelligence (AI). Here are just some of the capabilities:

- A Natural Language Processing (NLP) system that can understand conversations, allowing for a streamlined experience

- Computer vision software that verifies millions of images and documents like drivers' licenses and restaurant menus

- Sensor processing algorithms that help improve the accuracy in dense urban areas, including automatic crash detection by sensing unexpected movement from the phone of a driver or passenger

- Sophisticated machine learning algorithms that predict driver supply, rider demand, and ETAs

Such technologies are definitely amazing, but they are also required. There is no way that Uber could have scaled its growth—which has involved handling over 10 billion trips—without AI. In light of this, it should be no surprise that the company spends hundreds of millions on the technology and has a large group of AI experts on staff.[1]

But AI is not just for fast-charging startups. The technology is also proving a critical priority for traditional companies. Just look at McDonald's. In 2019, the company shelled out $300 million to acquire a tech startup, Dynamic Yield. It was the company's largest deal since it purchased Boston Market in 1999.[2]

---

[1]www.sec.gov/Archives/edgar/data/1543151/000119312519120759/d647752ds1a.htm#toc647752_11
[2]https://news.mcdonalds.com/news-releases/news-release-details/dynamic-yield-acquisition-release

Dynamic Yield, which was founded in 2011, is a pioneer in leveraging AI for creating personalized customer interactions across the Web, apps, and email. Some of its customers include the Hallmark Channel, IKEA, and Sephora.

As for McDonald's, it has been undergoing a digital transformation—and AI is a key part of the strategy. With Dynamic Yield, the company plans to use the technology to reimagine its Drive Thru, which accounts for a majority of its revenues. By analyzing data, such as the weather, traffic, and time of day, the digital menus will be dynamically changed to enhance the revenue opportunities. It also looks like McDonald's will use geofencing and even image recognition of license plates to enhance the targeting.

But this will just be the start. McDonald's expects to use AI for in-store kiosks and signage as well as the supply chain.

The company realizes that the future is both promising and dangerous. If companies are not proactive with new technologies, they may ultimately fail. Just look at how Kodak was slow to adapt to digital cameras. Or consider how the taxi industry did not change when faced with the onslaught of Uber and Lyft.

On the other hand, new technologies can be almost an elixir for a company. But there needs to be a solid strategy, a good understanding of what's possible, and a willingness to take risks. So in this book, I'll provide tools to help with all this.

OK then, how big will AI get? According to a study from PWC, it will add a staggering $15.7 trillion to the global GDP by 2030, which is more than the combined output of China and India. The authors of the report note: "AI touches almost every aspect of our lives. And it's only just getting started."[3]

True, when it comes to predicting trends, there can be a good deal of hype. However, AI may be different because it has the potential for turning into a general-purpose technology. A parallel to this is what happened in the nineteenth century with the emergence of electricity, which had a transformative impact across the world.

As a sign of the strategic importance of AI, tech companies like Google, Microsoft, Amazon.com, Apple, and Facebook have made substantial investments in this industry. For example, Google calls itself an "AI-first" company and has spent billions buying companies in the space as well as hiring thousands of data scientists.

In other words, more and more jobs will require knowledge of AI. Granted, this does not mean you'll need to learn programming languages or understand advanced statistics. But it will be critical to have a solid foundation of the fundamentals.

---

[3]www.pwc.com/gx/en/issues/data-and-analytics/publications/artificial-intelligence-study.html

As for this book, the goal is to provide actionable advice that can make a big difference in your organization and career. Now you will not find deeply technical explanations, code snippets, or equations. Instead, *Artificial Intelligence Basics* is about answering the top-of-mind questions that managers have: Where does AI make sense? What are the gotchas? How do you evaluate the technology? What about starting an AI pilot?

This book also takes a real-world view of the technology. A big advantage I have as a writer for Forbes.com and an advisor in the tech world is that I get to talk to many talented people in the AI field—and this helps me to identify what is really important in the industry. I also get to learn about case studies and examples of what works.

This book is organized in a way to cover the main topics in AI—and you do not have to read each chapter in order. *Artificial Intelligence Basics* is meant to be a handbook.

Here are brief descriptions of the chapters:

- *Chapter 1—AI Foundations*: This is an overview of the rich history of AI, which goes back to the 1950s. You will learn about brilliant researchers and computer scientists like Alan Turing, John McCarthy, Marvin Minsky, and Geoffrey Hinton. There will also be coverage of key concepts like the Turing Test, which gauges if a machine has achieved true AI.

- *Chapter 2—Data*: Data is the lifeblood of AI. It's how algorithms can find patterns and correlations to provide insights. But there are landmines with data, such as quality and bias. This chapter provides a framework to work with data in an AI project.

- *Chapter 3—Machine Learning*: This is a subset of AI and involves traditional statistical techniques like regressions. But in this chapter, we'll also cover the advanced algorithms, such as k-Nearest Neighbor (k-NN) and the Naive Bayes Classifier. Besides this, there will be a look at how to put together a machine learning model.

- *Chapter 4—Deep Learning*: This is another subset of AI and is clearly the one that has seen much of the innovation during the past decade. Deep learning is about using neural networks to find patterns that mimic the brain. In the chapter, we'll take a look at the main algorithms like recurrent neural networks (RNNs), convolutional neural networks (CNNs), and generative adversarial networks (GANs). There will also be explanations of key concepts like backpropagation.

- *Chapter 5—Robotic Process Automation:* This uses systems to automate repetitive processes, such as inputting data in a Customer Relationship Management (CRM) system. Robotic Process Automation (RPA) has seen tremendous growth during the past few years because of the high ROI (Return on Investment). The technology has also been an introductory way for companies to implement AI.

- *Chapter 6—Natural Language Processing (NLP)*: This form of AI, which involves understanding conversations, is the most ubiquitous as seen with Siri, Cortana, and Alexa. But NLP systems, such as chatbots, have also become critical in the corporate world. This chapter will show ways to use this technology effectively and how to avoid the tricky issues.

- *Chapter 7—Physical Robots*: AI is starting to have a major impact on this industry. With deep learning, it is getting easier for robots to understand their environments. In this chapter, we'll take a look at both consumer and industrial robots, such as with a myriad of use cases.

- *Chapter 8—Implementation of AI*: We'll take a step-by-step approach to putting together an AI project, from the initial concept to the deployment. This chapter will also cover the various tools like Python, TensorFlow, and PyTorch.

- *Chapter 9—The Future of AI:* This chapter will cover some of the biggest trends in AI like autonomous driving, weaponization of AI, technological unemployment, drug discovery, and regulation.

At the back of the book, you'll also find an appendix of resources for further study and a glossary of common terms related to AI.

## Accompanying Material

Any updates will be provided on my site at www.Taulli.com.

# AI Foundations

## History Lessons

*Artificial intelligence would be the ultimate version of Google. The ultimate search engine that would understand everything on the web. It would understand exactly what you wanted, and it would give you the right thing. We're nowhere near doing that now. However, we can get incrementally closer to that, and that is basically what we work on.*

—Larry Page, the co-founder of Google Inc. and
CEO of Alphabet[1]

In Fredric Brown's 1954 short story, "Answer," all of the computers across the 96 billion planets in the universe were connected into one super machine. It was then asked, "Is there a God?" to which it answered, "Yes, *now* there is a God."

No doubt, Brown's story was certainly clever—as well as a bit comical and chilling! Science fiction has been a way for us to understand the implications of new technologies, and artificial intelligence (AI) has been a major theme. Some of the most memorable characters in science fiction involve androids or computers that become self-aware, such as in *Terminator*, *Blade Runner*, *2001: A Space Odyssey,* and even *Frankenstein*.

But with the relentless pace of new technologies and innovation nowadays, science fiction is starting to become real. We can now talk to our smartphones and get answers; our social media accounts provide us with the content we're

---

[1] Founding CEO of Google Inc. The Academy of Achievement interview, www.achievement.org, October 28, 2000.

interested in; our banking apps provide us with reminders; and on and on. This personalized content creation almost seems magical but is quickly becoming normal in our everyday lives.

To understand AI, it's important to have a grounding in its rich history. You'll see how the development of this industry has been full of breakthroughs and setbacks. There is also a cast of brilliant researchers and academics, like Alan Turing, John McCarthy, Marvin Minsky, and Geoffrey Hinton, who pushed the boundaries of the technology. But through it all, there was constant progress.

Let's get started.

# Alan Turing and the Turing Test

Alan Turing is a towering figure in computer science and AI. He is often called the "father of AI."

In 1936, he wrote a paper called "On Computable Numbers." In it, he set forth the core concepts of a computer, which became known as the Turing machine. Keep in mind that real computers would not be developed until more than a decade later.

Yet it was his paper, called "Computing Machinery and Intelligence," that would become historic for AI. He focused on the concept of a machine that was intelligent. But in order to do this, there had to be a way to measure it. What is intelligence—at least for a machine?

This is where he came up with the famous "Turing Test." It is essentially a game with three players: two that are human and one that is a computer. The evaluator, a human, asks open-ended questions of the other two (one human, one computer) with the goal of determining which one is the human. If the evaluator cannot make a determination, then it is presumed that the computer is intelligent. Figure 1-1 shows the basic workflow of the Turing Test.



**Figure 1-1.** The basic workflow of the Turing Test

The genius of this concept is that there is no need to see if the machine actually knows something, is self-aware, or even if it is correct. Rather, the Turing Test indicates that a machine can process large amounts of information, interpret speech, and communicate with humans.

Turing believed that it would actually not be until about the turn of the century that a machine would pass his test. Yes, this was one of many predictions of AI that would come up short.

So how has the Turing Test held up over the years? Well, it has proven to be difficult to crack. Keep in mind that there are contests, such as the Loebner Prize and the Turing Test Competition, to encourage people to create intelligent software systems.

In 2014, there was a case where it did look like the Turing Test was passed. It involved a computer that said it was 13 years old.[2] Interestingly enough, the human judges likely were fooled because some of the answers had errors.

Then in May 2018 at Google's I/O conference, CEO Sundar Pichai gave a standout demo of Google Assistant.[3] Before a live audience, he used the device to call a local hairdresser to make an appointment. The person on the other end of the line acted as if she was talking to a person!

Amazing, right? Definitely. Yet it still probably did not pass the Turing Test. The reason is that the conversation was focused on one topic—not open ended.

As should be no surprise, there has been ongoing controversy with the Turing Test, as some people think it can be manipulated. In 1980, philosopher John Searle wrote a famous paper, entitled "Minds, Brains, and Programs," where he set up his own thought experiment, called the "Chinese room argument" to highlight the flaws.

Here's how it worked: Let's say John is in a room and does not understand the Chinese language. However, he does have manuals that provide easy-to-use rules to translate it. Outside the room is Jan, who does understand the language and submits characters to John. After some time, she will then get an accurate translation from John. As such, it's reasonable to assume that Jan believes that John can speak Chinese.

---

[2] www.theguardian.com/technology/2014/jun/08/super-computer-simulates-13-year-old-boy-passes-turing-test
[3] www.theverge.com/2018/5/8/17332070/google-assistant-makes-phone-call-demo-duplex-io-2018

Searle's conclusion:

> The point of the argument is this: if the man in the room does not understand Chinese on the basis of implementing the appropriate program for understanding Chinese then neither does any other digital computer solely on that basis because no computer, qua computer, has anything the man does not have.[4]

It was a pretty good argument—and has been a hot topic of debate in AI circles since.

Searle also believed there were two forms of AI:

- *Strong AI:* This is when a machine truly understands what is happening. There may even be emotions and creativity. For the most part, it is what we see in science fiction movies. This type of AI is also known as Artificial General Intelligence (AGI). Note that there are only a handful of companies that focus on this category, such as Google's DeepMind.

- *Weak AI:* With this, a machine is pattern matching and usually focused on narrow tasks. Examples of this include Apple's Siri and Amazon's Alexa.

The reality is that AI is in the early phases of weak AI. Reaching the point of strong AI could easily take decades. Some researchers think it may never happen.

Given the limitations to the Turing Test, there have emerged alternatives, such as the following:

- *Kurzweil-Kapor Test:* This is from futurologist Ray Kurzweil and tech entrepreneur Mitch Kapor. Their test requires that a computer carry on a conversation for two hours and that two of three judges believe it is a human talking. As for Kapor, he does not believe this will be achieved until 2029.

- *Coffee Test:* This is from Apple co-founder Steve Wozniak. According to the coffee test, a robot must be able to go into a stranger's home, locate the kitchen, and brew a cup of coffee.

---

[4] https://plato.stanford.edu/entries/chinese-room/

# The Brain Is a…Machine?

In 1943, Warren McCulloch and Walter Pitts met at the University of Chicago, and they became fast friends even though their backgrounds were starkly different as were their ages (McCulloch was 42 and Pitts was 18). McCulloch grew up in a wealthy Eastern Establishment family, having gone to prestigious schools. Pitts, on the other hand, grew up in a low-income neighborhood and was even homeless as a teenager.

Despite all this, the partnership would turn into one of the most consequential in the development of AI. McCulloch and Pitts developed new theories to explain the brain, which often went against the conventional wisdom of Freudian psychology. But both of them thought that logic could explain the power of the brain and also looked at the insights from Alan Turing. From this, they co-wrote a paper in 1943 called "A Logical Calculus of the Ideas Immanent in Nervous Activity," and it appeared in the *Bulletin of Mathematical Biophysics*. The thesis was that the brain's core functions like neurons and synapses could be explained by logic and mathematics, say with logical operators like And, Or, and Not. With these, you could construct a complex network that could process information, learn, and think.

Ironically, the paper did not get much traction with neurologists. But it did get the attention with those working on computers and AI.

# Cybernetics

While Norbert Wiener created various theories, his most famous one was about cybernetics. It was focused on understanding control and communications with animals, people, and machines—showing the importance of feedback loops.

In 1948, Wiener published *Cybernetics: Or Control and Communication in the Animal and the Machine*. Even though it was a scholarly work—filled with complex equations—the book still became a bestseller, hitting the *New York Times* list.

It was definitely wide ranging. Some of the topics included Newtonian mechanics, meteorology, statistics, astronomy, and thermodynamics. This book would anticipate the development of chaos theory, digital communications, and even computer memory.

But the book would also be influential for AI. Like McCulloch and Pitts, Wiener compared the human brain to the computer. Furthermore, he speculated that a computer would be able to play chess and eventually beat grand masters. The main reason is that he believed that a machine could learn as it played games. He even thought that computers would be able to replicate themselves.

But *Cybernetics* was not utopian either. Wiener was also prescient in understanding the downsides of computers, such as the potential for dehumanization. He even thought that machines would make people unnecessary.

It was definitely a mixed message. But Wiener's ideas were powerful and spurred the development of AI.

# The Origin Story

John McCarthy's interest in computers was spurred in 1948, when he attended a seminar, called "Cerebral Mechanisms in Behavior," which covered the topic of how machines would eventually be able to think. Some of the participants included the leading pioneers in the field such as John von Neumann, Alan Turing, and Claude Shannon.

McCarthy continued to immerse himself in the emerging computer industry—including a stint at Bell Labs—and in 1956, he organized a ten-week research project at Dartmouth University. He called it a "study of artificial intelligence." It was the first time the term had been used.

The attendees included academics like Marvin Minsky, Nathaniel Rochester, Allen Newell, O. G. Selfridge, Raymond Solomonoff, and Claude Shannon. All of them would go on to become major players in AI.

The goals for the study were definitely ambitious:

> The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for a summer.[5]

At the conference, Allen Newell, Cliff Shaw, and Herbert Simon demoed a computer program called the Logic Theorist, which they developed at the Research and Development (RAND) Corporation. The main inspiration came from Simon (who would win the Nobel Prize in Economics in 1978). When he saw how computers printed out words on a map for air defense systems, he realized that these machines could be more than just about processing numbers. It could also help with images, characters, and symbols—all of which could lead to a thinking machine.

Regarding Logic Theorist, the focus was on solving various math theorems from *Principia Mathematica*. One of the solutions from the software turned out to be more elegant—and the co-author of the book, Bertrand Russell, was delighted.

---

[5] www-formal.stanford.edu/jmc/history/dartmouth/dartmouth.html

Creating the Logic Theorist was no easy feat. Newell, Shaw, and Simon used an IBM 701, which used machine language. So they created a high-level language, called IPL (Information Processing Language), that sped up the programming. For several years, it was the language of choice for AI.

The IBM 701 also did not have enough memory for the Logic Theorist. This led to another innovation: list processing. It allowed for dynamically allocating and deallocating memory as the program ran.

Bottom line: The Logic Theorist is considered the first AI program ever developed.

Despite this, it did not garner much interest! The Dartmouth conference was mostly a disappointment. Even the phrase "artificial intelligence" was criticized.

Researchers tried to come up with alternatives, such as "complex information processing." But they were not catchy like AI was—and the term stuck.

As for McCarthy, he continued on his mission to push innovation in AI. Consider the following:

- During the late 1950s, he developed the Lisp programming language, which was often used for AI projects because of the ease of using nonnumerical data. He also created programming concepts like recursion, dynamic typing, and garbage collection. Lisp continues to be used today, such as with robotics and business applications. While McCarthy was developing the language, he also co-founded the MIT Artificial Intelligence Laboratory.

- In 1961, he formulated the concept of time-sharing of computers, which had a transformative impact on the industry. This also led to the development of the Internet and cloud computing.

- A few years later, he founded Stanford's Artificial Intelligence Laboratory.

- In 1969, he wrote a paper called "Computer-Controlled Cars," in which he described how a person could enter directions with a keyboard and a television camera would navigate the vehicle.

- He won the Turing Award in 1971. This prize is considered the Nobel Prize for Computer Science.

In a speech in 2006, McCarthy noted that he was too optimistic about the progress of strong AI. According to him, "we humans are not very good at identifying the heuristics we ourselves use."[6]

# Golden Age of AI

From 1956 to 1974, the AI field was one of the hottest spots in the tech world. A major catalyst was the rapid development in computer technologies. They went from being massive systems—based on vacuum tubes—to smaller systems run on integrated circuits that were much quicker and had more storage capacity.

The federal government was also investing heavily in new technologies. Part of this was due to the ambitious goals of the Apollo space program and the heavy demands of the Cold War.

As for AI, the main funding source was the Advanced Research Projects Agency (ARPA), which was launched in the late 1950s after the shock of Russia's Sputnik. The spending on projects usually came with few requirements. The goal was to inspire breakthrough innovation. One of the leaders of ARPA, J. C. R. Licklider, had a motto of "fund people, not projects." For the most part, the majority of the funding was from Stanford, MIT, Lincoln Laboratories, and Carnegie Mellon University.

Other than IBM, the private sector had little involvement in AI development. Keep in mind that—by the mid-1950s—IBM would pull back and focus on the commercialization of its computers. There was actually fear from customers that this technology would lead to significant job losses. So IBM did not want to be blamed.

In other words, much of the innovation in AI spun out from academia. For example, in 1959, Newell, Shaw, and Simon continued to push the boundaries in the AI field with the development of a program called "General Problem Solver." As the name implied, it was about solving math problems, such as the Tower of Hanoi.

But there were many other programs that attempted to achieve some level of strong AI. Examples included the following:

- *SAINT or Symbolic Automatic INTegrator (1961)*: This program, created by MIT researcher James Slagle, helped to solve freshman calculus problems. It would be updated into other programs, called SIN and MACSYMA, that did much more advanced math. SAINT was actually the first example of an expert system, a category of AI we'll cover later in this chapter.

---

[6] www.technologyreview.com/s/425913/computing-pioneer-dies/

- *ANALOGY (1963)*: This program was the creation of MIT professor Thomas Evans. The application demonstrated that a computer could solve analogy problems of an IQ test.

- *STUDENT (1964)*: Under the supervision of Minsky at MIT, Daniel Bobrow created this AI application for his PhD thesis. The system used Natural Language Processing (NLP) to solve algebra problems for high school students.

- *ELIZA (1965)*: MIT professor Joseph Weizenbaum designed this program, which instantly became a big hit. It even got buzz in the mainstream press. It was named after Eliza (based on George Bernard Shaw's play *Pygmalion*) and served as a psychoanalyst. A user could type in questions, and ELIZA would provide counsel (this was the first example of a chatbot). Some people who used it thought the program was a real person, which deeply concerned Weizenbaum since the underlying technology was fairly basic. You can find examples of ELIZA on the web, such as at http://psych.fullerton.edu/mbirnbaum/psych101/Eliza.htm.

- *Computer Vision (1966)*: In a legendary story, MIT's Marvin Minsky said to a student, Gerald Jay Sussman, to spend the summer linking a camera to a computer and getting the computer to describe what it saw. He did just that and built a system that detected basic patterns. It was the first use of computer vision.

- *Mac Hack (1968)*: MIT professor Richard D. Greenblatt created this program that played chess. It was the first to play in real tournaments and got a C-rating.

- *Hearsay I (Late 1960s)*: Professor Raj Reddy developed a continuous speech recognition system. Some of his students would then go on to create Dragon Systems, which became a major tech company.

During this period, there was a proliferation of AI academic papers and books. Some of the topics included Bayesian methods, machine learning, and vision.

But there were generally two major theories about AI. One was led by Minsky, who said that there needed to be symbolic systems. This meant that AI should be based on traditional computer logic or preprogramming—that is, the use of approaches like If-Then-Else statements.

Next, there was Frank Rosenblatt, who believed that AI needed to use systems similar to the brain like neural networks (this field was also known as connectionism). But instead of calling the inner workings neurons, he referred to them as perceptrons. A system would be able to learn as it ingested data over time.

In 1957, Rosenblatt created the first computer program for this, called the Mark 1 Perceptron. It included cameras to help to differentiate between two images (they had 20 × 20 pixels). The Mark 1 Perceptron would use data that had random weightings and then go through the following process:

1.  Take in an input and come up with the perceptron output.

2.  If there is not a match, then

    a.  If the output should have been 0 but was 1, then the weight for 1 will be decreased.

    b.  If the output should have been 1 but was 0, then the weight for 1 will be increased.

3.  Repeat steps #1 and #2 until the results are accurate.

This was definitely pathbreaking for AI. The *New York Times* even had a write-up for Rosenblatt, extolling "The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."[7]

But there were still nagging issues with the perceptron. One was that the neural network had only one layer (primarily because of the lack of computation power at the time). Next, brain research was still in the nascent stages and did not offer much in terms of understanding cognitive ability.

Minsky would co-write a book, along with Seymour Papert, called *Perceptrons* (1969). The authors were relentless in attacking Rosenblatt's approach, and it quickly faded away. Note that in the early 1950s Minsky developed a crude neural net machine, such as by using hundreds of vacuum tubes and spare parts from a B-24 bomber. But he saw that the technology was nowhere at a point to be workable.

Rosenblatt tried to fight back, but it was too late. The AI community quickly turned sour on neural networks. Rosenblatt would then die a couple years later in a boating accident. He was 43 years old.

Yet by the 1980s, his ideas would be revived—which would lead to a revolution in AI, primarily with the development of deep learning.

---

[7] www.nytimes.com/1958/07/08/archives/new-navy-device-learns-by-doing-psychologist-shows-embryo-of.html

For the most part, the Golden Age of AI was freewheeling and exciting. Some of the brightest academics in the world were trying to create machines that could truly think. But the optimism often went to the extremes. In 1965, Simon said that within 20 years, a machine could do anything a human could. Then in 1970, in an interview with *Life* magazine, he said this would happen in only 3–8 years (by the way, he was an advisor on the *2001: A Space Odyssey* movie).

Unfortunately, the next phase of AI would be much darker. There were more academics who were becoming skeptical. Perhaps the most vocal was Hubert Dreyfus, a philosopher. In books such as *What Computers Still Can't Do: A Critique of Artificial Reason,*[8] he set forth his ideas that computers were not similar to the human brain and that AI would woefully fall short of the lofty expectations.

# AI Winter

During the early 1970s, the enthusiasm for AI started to wane. It would become known as the "AI winter," which would last through 1980 or so (the term came from "nuclear winter," an extinction event where the sun is blocked and temperatures plunge across the world).

Even though there were many strides made with AI, they still were mostly academic and involved in controlled environments. At the time, the computer systems were still limited. For example, a DEC PDP-11/45—which was common for AI research—had the ability to expand its memory to only 128K.

The Lisp language also was not ideal for computer systems. Rather, in the corporate world, the focus was primarily on FORTRAN.

Next, there were still many complex aspects when understanding intelligence and reasoning. Just one is disambiguation. This is the situation when a word has more than one meaning. This adds to the difficulty for an AI program since it will also need to understand the context.

Finally, the economic environment in the 1970s was far from robust. There were persistent inflation, slow growth, and supply disruptions, such as with the oil crisis.

Given all this, it should be no surprise that the US government was getting more stringent with funding. After all, for a Pentagon planner, how useful is a program that can play chess, solve a theorem, or recognize some basic images?

Not much, unfortunately.

---

[8] MIT Press, 1972.

A notable case is the Speech Understanding Research program at Carnegie Mellon University. For the Defense Advanced Research Projects Agency (DARPA), it thought this speech recognition system could be used for fighter pilots to make voice commands. But it proved to be unworkable. One of the programs, which was called Harpy, could understand 1,011 words—which is what a typical 3-year-old knows.

The officials at DARPA actually thought that it had been hoodwinked and eliminated the $3 million annual budget for the program.

But the biggest hit to AI came via a report—which came out in 1973—from Professor Sir James Lighthill. Funded by the UK Parliament, it was a full-on repudiation of the "grandiose objectives" of strong AI. A major issue he noted was "combinatorial explosion," which was the problem where the models got too complicated and were difficult to adjust.

The report concluded: "In no part of the field have the discoveries made so far produced the major impact that was then promised."[9] He was so pessimistic that he did not believe computers would be able to recognize images or beat a chess grand master.

The report also led to a public debate that was televised on the BCC (you can find the videos on YouTube). It was Lighthill against Donald Michie, Richard Gregory, and John McCarthy.

Even though Lighthill had valid points—and evaluated large amounts of research—he did not see the power of weak AI. But it did not seem to matter as the winter took hold.

Things got so bad that many researchers changed their career paths. And as for those who still studied AI, they often referred to their work with other terms—like machine learning, pattern recognition, and informatics!

# The Rise and Fall of Expert Systems

Even during the AI winter, there continued to be major innovations. One was backpropagation, which is essential for assigning weights for neural networks. Then there was the development of the recurrent neural network (RNN). This allows for connections to move through the input and output layers.

But in the 1980s and 1990s, there also was the emergence of expert systems. A key driver was the explosive growth of PCs and minicomputers.

---

[9] The 1973 "Artificial Intelligence: A General Survey" by Professor Sir James Lighthill of Cambridge University, www.bbc.com/timelines/zq376fr.

Expert systems were based on the concepts of Minsky's symbolic logic, involving complex pathways. These were often developed by domain experts in particular fields like medicine, finance, and auto manufacturing.

Figure 1-2 shows the key parts of an expert system.



**Figure 1-2.** Key parts of an expert system

While there are expert systems that go back to the mid-1960s, they did not gain commercial use until the 1980s. An example was XCON (eXpert CONfigurer), which John McDermott developed at Carnegie Mellon University. The system allowed for optimizing the selection of computer components and initially had about 2,500 rules. Think of it as the first recommendation engine. From the launch in 1980, it turned out to be a big cost saver for DEC for its line of VAX computers (about $40 million by 1986).

When companies saw the success of XCON, there was a boom in expert systems—turning into a billion-dollar industry. The Japanese government also saw the opportunity and invested hundreds of millions to bolster its home market. However, the results were mostly a disappointment. Much of the innovation was in the United States.

Consider that IBM used an expert system for its Deep Blue computer. In 1996, it would beat grand chess master Garry Kasparov, in one of six matches. Deep Blue, which IBM had been developing since 1985, processed 200 million positions per second.

But there were issues with expert systems. They were often narrow and difficult to apply across other categories. Furthermore, as the expert systems got larger, it became more challenging to manage them and feed data. The

result was that there were more errors in the outcomes. Next, testing the systems often proved to be a complex process. Let's face it, there were times when the experts would disagree on fundamental matters. Finally, expert systems did not learn over time. Instead, there had to be constant updates to the underlying logic models, which added greatly to the costs and complexities.

By the late 1980s, expert systems started to lose favor in the business world, and many startups merged or went bust. Actually, this helped cause another AI winter, which would last until about 1993. PCs were rapidly eating into higher-end hardware markets, which meant a steep reduction in Lisp-based machines.

Government funding for AI, such as from DARPA, also dried up. Then again, the Cold War was rapidly coming to a quiet end with the fall of the Soviet Union.

# Neural Networks and Deep Learning

As a teen in the 1950s, Geoffrey Hinton wanted to be a professor and to study AI. He came from a family of noted academics (his great-great-grandfather was George Boole). His mom would often say, "Be an academic or be a failure."[10]

Even during the first AI winter, Hinton was passionate about AI and was convinced that Rosenblatt's neural network approach was the right path. So in 1972, he received his PhD on the topic from the University of Edinburgh.

But during this period, many people thought that Hinton was wasting his time and talents. AI was essentially considered a fringe area. It wasn't even thought of as a science.

But this only encouraged Hinton more. He relished his position as an outsider and knew that his ideas would win out in the end.

Hinton realized that the biggest hindrance to AI was computer power. But he also saw that time was on his side. Moore's Law predicted that the number of components on a chip would double about every 18 months.

In the meantime, Hinton worked tirelessly on developing the core theories of neural networks—something that eventually became known as deep learning. In 1986, he wrote—along with David Rumelhart and Ronald J. Williams—a pathbreaking paper, called "Learning Representations by Back-propagating Errors." It set forth key processes for using backpropagation in neural networks. The result was that there would be significant improvement in accuracy, such as with predictions and visual recognition.

---

[10] https://torontolife.com/tech/ai-superstars-google-facebook-apple-studied-guy/

Of course, this did not happen in isolation. Hinton's pioneering work was based on the achievements of other researchers who also were believers in neural networks. And his own research spurred a flurry of other major achievements:

- *1980*: Kunihiko Fukushima created Neocognitron, which was a system to recognize patterns that became the basis of convolutional neural networks. These were based on the visual cortex of animals.

- *1982*: John Hopfield developed "Hopfield Networks." This was essentially a recurrent neural network.

- *1989*: Yann LeCun merged convolutional networks with backpropagation. This approach would find use cases with analyzing handwritten checks.

- *1989*: Christopher Watkins' PhD thesis, "Learning from Delayed Rewards," described Q-Learning. This was a major advance in helping with reinforcement learning.

- *1998*: Yann LeCun published "Gradient-Based Learning Applied to Document Recognition," which used descent algorithms to improve neural networks.

# Technological Drivers of Modern AI

Besides advances in new conceptual approaches, theories, and models, AI had some other important drivers. Here's a look at the main ones:

- *Explosive Growth in Datasets*: The internet has been a major factor for AI because it has allowed for the creation of massive datasets. In the next chapter, we'll take a look at how data has transformed this technology.

- *Infrastructure*: Perhaps the most consequential company for AI during the past 15 years or so has been Google. To keep up with the indexing of the Web—which was growing at a staggering rate—the company had to come up with creative approaches to build scalable systems. The result has been innovation in commodity server clusters, virtualization, and open source software. Google was also one of the early adopters of deep learning, with the launch of the "Google Brain" project in 2011. Oh, and a few years later the company hired Hinton.

- *GPUs (Graphics Processing Units)*: This chip technology, which was pioneered by NVIDIA, was originally for high-speed graphics in games. But the architecture of GPUs would eventually be ideal for AI as well. Note that most deep learning research is done with these chips. The reason is that—with parallel processing—the speed is multiples higher than traditional CPUs. This means that computing a model may take a day or two vs. weeks or even months.

All these factors reinforced themselves—adding fuel to the growth of AI. What's more, these factors are likely to remain vibrant for many years to come.

## Structure of AI

In this chapter, we've covered many concepts. Now it can be tough to understand the organization of AI. For instance, it is common to see terms like machine learning and deep learning get confused. But it is essential to understand the distinctions, which we will cover in detail in the rest of this book.

But on a high-level view of things, Figure 1-3 shows how the main elements of AI relate to each other. At the top is AI, which covers a wide variety of theories and technologies. You can then break this down into two main categories: machine learning and deep learning.



**Figure 1-3.** This is a high-level look at the main components of the AI world

# Conclusion

There's nothing new that AI is a buzzword today. The term has seen various stomach-churning boom-bust cycles.

Maybe it will once again go out of favor? Perhaps. But this time around, there are true innovations with AI that are transforming businesses. Mega tech companies like Google, Microsoft, and Facebook consider the category to be a major priority. All in all, it seems like a good bet that AI will continue to grow and change our world.

# Key Takeaways

- Technology often takes longer to evolve than originally understood.

- AI is not just about computer science and mathematics. There have been key contributions from fields like economics, neuroscience, psychology, linguistics, electrical engineering, mathematics, and philosophy.

- There are two main types of AI: weak and strong. Strong is where machines become self-aware, whereas weak is for systems that focus on specific tasks. Currently, AI is at the weak stage.

- The Turing Test is a common way to test if a machine can think. It is based on whether someone really thinks a system is intelligent.

- Some of the key drivers for AI include new theories from researchers like Hinton, the explosive growth in data, new technology infrastructure, and GPUs.

# Data

## The Fuel for AI

Pinterest is one of the hottest startups in Silicon Valley, allowing users to pin their favorite items to create engaging boards. The site has 250 million MAUs (monthly active users) and posted $756 million in revenue in 2018.[1]

A popular activity for Pinterest is to plan for weddings. The bride-to-be will have pins for gowns, venues, honeymoon spots, cakes, invitations, and so on.

This also means that Pinterest has the advantage of collecting huge amounts of valuable data. Part of this helps provide for targeted ads. Yet there are also opportunities for email campaigns. In one case, Pinterest sent one that said:

> You're getting married! And because we love wedding planning—especially all the lovely stationery—we invite you to browse our best boards curated by graphic designers, photographers and fellow brides-to-be, all Pinners with a keen eye and marriage on the mind.[2]

The problem: Plenty of the recipients of the email were already married or not expecting to marry anytime soon.

---

[1] www.cnbc.com/2019/03/22/pinterest-releases-s-1-for-ipo.html
[2] www.businessinsider.com/pinterest-accidental-marriage-emails-2014-9

Pinterest did act quickly and put out this apology:

> Every week, we email collections of category-specific pins and boards to pinners we hope will be interested in them. Unfortunately, one of these recent emails suggested that pinners were actually getting married, rather than just potentially interested in wedding-related content. We're sorry we came off like an overbearing mother who is always asking when you'll find a nice boy or girl.

It's an important lesson. Even some of the most tech-savvy companies blow it.

For example, there are some cases where the data may be spot-on but the outcome could still be an epic failure. Consider the case with Target. The company leveraged its massive data to send personalized offers to expectant mothers. This was based on those customers who made certain types of purchases, such as for unscented lotions. Target's system would create a pregnancy score that even provided estimates of due dates.

Well, the father of one of the customers saw the email and was furious, saying his daughter was not pregnant.[3]

But she was—and yes, she had been hiding this fact from her father.

There's no doubt that data is extremely powerful and critical for AI. But you need to be thoughtful and understand the risks. In this chapter, we'll take a look at some of the things you need to know.

# Data Basics

It's good to have an understanding of the jargon of data.

First of all, a bit (which is short for "binary digit") is the smallest form of data in a computer. Think of it as the atom. A bit can either be 0 or 1, which is binary. It is also generally used to measure the amount of data that is being transferred (say within a network or the Internet).

A byte, on the other hand, is mostly for storage. Of course, the numbers of bytes can get large very fast. Let's see how in Table 2-1.

---

[3] www.businessinsider.com/the-incredible-story-of-how-target-exposed-a-teen-girls-pregnancy-2012-2

**Table 2-1.** Types of data levels

| Unit | Value | Use Case |
|------|-------|----------|
| Megabyte | 1,000 kilobytes | A small book |
| Gigabyte | 1,000 megabytes | About 230 songs |
| Terabyte | 1,000 gigabytes | 500 hours of movies |
| Petabyte | 1,000 terabytes | Five years of the Earth Observing System (EOS) |
| Exabyte | 1,000 petabytes | The entire Library of Congress 3,000 times over |
| Zettabyte | 1,000 exabytes | 36,000 years of HD-TV video |
| Yottabytes | 1,000 zettabytes | This would require a data center the size of Delaware and Rhode Island combined |

Data can also come from many different sources. Here is just a sampling:

- Web/social (Facebook, Twitter, Instagram, YouTube)
- Biometric data (fitness trackers, genetics tests)
- Point of sale systems (from brick-and-mortar stores and e-commerce sites)
- Internet of Things or IoT (ID tags and smart devices)
- Cloud systems (business applications like Salesforce.com)
- Corporate databases and spreadsheets

# Types of Data

There are four ways to organize data. First, there is structured data, which is usually stored in a relational database or spreadsheet. Some examples include the following:

- Financial information
- Social Security numbers
- Addresses
- Product information
- Point of sale data
- Phone numbers

For the most part, structured data is easier to work with. This data often comes from CRM (Customer Relationship Management) and ERP (Enterprise Resource Planning) systems—and usually has lower volumes. It also tends to

be more straightforward, say in terms of analysis. There are various BI (Business Intelligence) programs that can help derive insights from structured data. However, this type of data accounts for about 20% of an AI project.

The majority will instead come from unstructured data, which is information that has no predefined formatting. You'll have to do this yourself, which can be tedious and time consuming. But there are tools like next-generation databases—such as those based on NoSQL—that can help with the process. AI systems are also effective in terms of managing and structuring the data, as the algorithms can recognize patterns.

Here are examples of unstructured data:

- Images
- Videos
- Audio files
- Text files
- Social network information like tweets and posts
- Satellite images

Now there is some data that is a hybrid of structured and unstructured sources—called semi-structured data. The information has some internal tags that help with categorization.

Examples of semi-structured data include XML (Extensible Markup Language), which is based on various rules to identify elements of a document, and JSON (JavaScript Object Notation), which is a way to transfer information on the Web through APIs (Application Programming Interfaces).

But semi-structured data represents only about 5% to 10% of all data.

Finally, there is time-series data, which can be both for structured, unstructured, and semi-structured data. This type of information is for interactions, say for tracking the "customer journey." This would be collecting information when a user goes to the web site, uses an app, or even walks into a store.

Yet this kind of data is often messy and difficult to understand. Part of this is due to understanding the intent of the users, which can vary widely. There is also huge volumes of interactional data, which can involve trillions of data points. Oh, and the metrics for success may not be clear. Why is a user doing something on the site?

But AI is likely to be critical for such issues. Although, for the most part, the analysis of time-series data is still in the early stages.

# Big Data

With the ubiquity of Internet access, mobile devices, and wearables, there has been the unleashing of a torrent of data. Every second, Google processes over 40,000 searches or 3.5 billion a day. On a minute-by-minute basis, Snapchat users share 527,760 photos, and YouTube users watch more than 4.1 million videos. Then there are the old-fashioned systems, like emails, that continue to see significant growth. Every minute, there are 156 million messages sent.[4]

But there is something else to consider: Companies and machines also generate huge sums of data. According to research from Statista, the number of sensors will reach 12.86 billion by 2020.[5]

In light of all this, it seems like a good bet that the volumes of data will continue to increase at a rapid clip. In a report from International Data Corporation (IDC) called "Data Age 2025," the amount of data created is expected to hit a staggering 163 zettabytes by 2025.[6] This is about ten times the amount in 2017.

To deal with all this, there has emerged a category of technology called Big Data. This is how Oracle explains the importance of this trend:

> Today, big data has become capital. Think of some of the world's biggest tech companies. A large part of the value they offer comes from their data, which they're constantly analyzing to produce more efficiency and develop new products.[7]

So yes, Big Data will remain a critical part of many AI projects.

Then what exactly is Big Data? What's a good definition? Actually, there isn't one, even though there are many companies that focus on this market! But Big Data does have the following characteristics, which are called the three Vs (Gartner analyst Doug Laney came up with this structure back in 2001[8]): volume, variety, and velocity.

---

[4] www.forbes.com/sites/bernardmarr/2018/05/21/how-much-data-do-we-create-every-day-the-mind-blowing-stats-everyone-should-read/#788c13c660ba

[5] www.forbes.com/sites/louiscolumbus/2018/06/06/10-charts-that-will-challenge-your-perspective-of-iots-growth/#4e9fac23ecce

[6] https://blog.seagate.com/business/enormous-growth-in-data-is-coming-how-to-prepare-for-it-and-prosper-from-it/

[7] www.oracle.com/big-data/guide/what-is-big-data.html

[8] https://blogs.gartner.com/doug-laney/files/2012/01/ad949-3D-Data-Management-Controlling-Data-Volume-Velocity-and-Variety.pdf

## Volume

This is the scale of the data, which is often unstructured. There is no hard-and-fast rule on a threshold, but it is usually tens of terabytes.

Volume is often a major challenge when it comes to Big Data. But cloud computing and next-generation databases have been a big help—in terms of capacity and lower costs.

## Variety

This describes the diversity of the data, say a combination of structured, semi-structured, and unstructured data (explained above). It also shows the different sources of the data and uses. No doubt, the high growth in unstructured data has been a key to the variety of Big Data.

Managing this can quickly become a major challenge. Yet machine learning is often something that can help streamline the process.

## Velocity

This shows the speed at which data is being created. As seen earlier in this chapter, services like YouTube and Snapchat have extreme levels of velocity (this is often referred to as a "firehouse" of data). This requires heavy investments in next-generation technologies and data centers. The data is also often processed in memory not with disk-based systems.

Because of these issues, velocity is often considered the most difficult when it comes to the three Vs. Let's face it, in today's digital world, people want their data as fast as possible. If it is too slow, people will get frustrated and go somewhere else.

Over the years, though, as Big Data has evolved, there have been more Vs added. Currently, there are over ten.

But here are some of the common ones:

- *Veracity*: This is about data that is deemed accurate. In this chapter, we'll look at some of the techniques to evaluate veracity.

- *Value*: This shows the usefulness of the data. Often this is about having a trusted source.

- *Variability*: This means that data will usually change over time. For example, this is the case with social media content that can morph based on overall sentiment regarding new developments and breaking news.

- *Visualization*: This is using visuals—like graphs—to better understand the data.

As you can see, managing Big Data has many moving parts, which leads to complexity. This helps to explain why many companies still use only a tiny fraction of their data.

# Databases and Other Tools

There are a myriad of tools that help with data. At the core of this is the database. As should be no surprise, there has been an evolution of this critical technology over the decades. But even older technologies like relational databases are still very much in use today. When it comes to mission-critical data, companies are reluctant to make changes—even if there are clear benefits.

To understand this market, let's rewind back to 1970, when IBM computer scientist Edgar Codd published "A Relational Model of Data for Large Shared Data Banks." It was pathbreaking as it introduced the structure of relational databases. Up until this point, databases were fairly complex and rigid—structured as hierarchies. This made it time consuming to search and find relationships in the data.

As for Codd's relational database approach, it was built for more modern machines. The SQL script language was easy to use allowing for CRUD (Create, Read, Update, and Delete) operations. Tables also had connections with primary and foreign keys, which made important connections like the following:

- *One-to-One*: One row in a table is linked to only one row in another table. Example: A driver's license number, which is unique, is associated with one employee.

- *One-to-Many*: This is where one row in a table is linked to other tables. Example: A customer has multiple purchase orders.

- *Many-to-Many*: Rows from one table are associated with rows of another. Example: Various reports have various authors.

With these types of structures, a relational database could streamline the process of creating sophisticated reports. It truly was revolutionary.

But despite the advantages, IBM was not interested in the technology and continued to focus on its proprietary systems. The company thought that the relational databases were too slow and brittle for enterprise customers.

But there was someone who had a different opinion on the matter: Larry Ellison. He read Codd's paper and knew it would be a game changer. To prove this, he would go on to co-found Oracle in 1977 with a focus on building relational databases—which would quickly become a massive market. Codd's paper was essentially a product roadmap for his entrepreneurial efforts.

It was not until 1993 that IBM came out with its own relational database, DB2. But it was too late. By this time, Oracle was the leader in the database market.

Through the 1980s and 1990s, the relational database was the standard for mainframe and client-server systems. But when Big Data became a factor, the technology had serious flaws like the following:

- *Data Sprawl*: Over time, different databases would spread across an organization. The result was that it got tougher to centralize the data.

- *New Environments*: Relational database technology was not built for cloud computing, high-velocity data, or unstructured data.

- *High Costs*: Relational databases can be expensive. This means that it can be prohibitive to use the technology for AI projects.

- *Development Challenges*: Modern software development relies heavily on iterating. But relational databases have proven challenging for this process.

In the late 1990s, there were open source projects developed to help create next-generation database systems. Perhaps the most critical one came from Doug Cutting who developed Lucene, which was for text searching. The technology was based on a sophisticated index system that allowed for low-latency performance. Lucene was an instant hit, and it started to evolve, such as with Apache Nutch that efficiently crawled the Web and stored the data in an index.

But there was a big problem: To crawl the Web, there needed to be an infrastructure that could hyperscale. So in late 2003, Cutting began development on a new kind of infrastructure platform that could solve the problem. He got the idea from a paper published from Google, which described its massive file system. A year later, Cutting had built his new platform, which allowed for sophisticated storage without the complexity. At the core of this was MapReduce that allowed processing across multiple servers. The results would then be merged, allowing for meaningful reports.

Eventually, Cutting's system morphed into a platform called Hadoop—and it would be essential for managing Big Data, such as making it possible to create sophisticated data warehouses. Initially, Yahoo! used it, and then it quickly spread, as companies like Facebook and Twitter adopted the technology. These companies were now able to get a 360 view of their data, not just subsets. This meant there could be more effective data experiments.

But as an open source project, Hadoop still lacked the sophisticated systems for enterprise customers. To deal with this, a startup called Hortonworks built new technologies like YARN on top of the Hadoop platform. It had features like in-memory analytic processing, online data processing, and interactive SQL processing. Such capabilities supported adoption of Hadoop across many corporations.

But of course, there emerged other open source data warehouse projects. The well-known ones, like Storm and Spark, focused on streaming data. Hadoop, on the other hand, was optimized for batch processing.

Besides data warehouses, there was also innovation of the traditional database business. Often these were known as NoSQL systems. Take MongoDB. It started as an open source project and has turned into a highly successful company, which went public in October 2017. The MongoDB database, which has over 40 million downloads, is built to handle cloud, on-premise, and hybrid environments.[9] There is also much flexibility structuring the data, which is based on a document model. MongoDB can even manage structured and unstructured data at high petabyte scale.

Even though startups have been a source of innovation in database systems and storage, it's important to note that the mega tech operators have also been critical. Then again, companies like Amazon.com and Google have had to find ways to deal with the huge scale of data because of the need for managing their massive platforms.

One of the innovations has been the data lake, which allows for seamless storage of structured and unstructured data. Note that there is no need to reformat the data. A data lake will handle this and allow you to quickly perform AI functions. According to a study from Aberdeen, companies who use this technology have an average of 9% organic growth compared to those who do not.[10]

Now this does not mean you have to get rid of your data warehouses. Rather, both serve particular functions and use cases. A data warehouse is generally good for structured data, whereas a data lake is better for diverse environments. What's more, it's likely that a large portion of the data will never be used.

---

[9] www.mongodb.com/what-is-mongodb
[10] https://aws.amazon.com/big-data/datalakes-and-analytics/what-is-a-data-lake/

For the most part, there are a myriad of tools. And expect more to be developed as data environments get more complex.

But this does not mean you should chose the latest technology. Again, even older relational databases can be quite effective with AI projects. The key is understanding the pros/cons of each and then putting together a clear strategy.

# Data Process

The amount of money shelled out on data is enormous. According to IDC, the spending on Big Data and analytics solutions is forecasted to go from $166 billion in 2018 to $260 billion by 2022.[11] This represents an 11.9% compound annual growth rate. The biggest spenders include banks, discrete manufacturers, process manufacturers, professional service firms, and the federal government. They account for close to half the overall amount.

Here's what IDC's Jessica Goepfert—the program vice president (VP) of Customer Insights and Analysis—said:

> *At a high level, organizations are turning to Big Data and analytics solutions to navigate the convergence of their physical and digital worlds. This transformation takes a different shape depending on the industry. For instance, within banking and retail—two of the fastest growth areas for Big Data and analytics—investments are all about managing and reinvigorating the customer experience. Whereas in manufacturing, firms are reinventing themselves to essentially be high tech companies, using their products as a platform to enable and deliver digital services.[12]*

But a high level of spending does not necessarily translate into good results. A Gartner study estimates that roughly 85% of Big Data projects are abandoned before they get to the pilot stage.[13] Some of the reasons include the following:

- Lack of a clear focus

- Dirty data

- Investment in the wrong IT tools

- Problems with data collection

- Lack of buy-in from key stakeholders and champions in the organization

---

[11] www.idc.com/getdoc.jsp?containerId=prUS44215218
[12] www.idc.com/getdoc.jsp?containerId=prUS44215218
[13] www.techrepublic.com/article/85-of-big-data-projects-fail-but-your-developers-can-help-yours-succeed/

In light of this, it is critical to have a data process. Notwithstanding there are many approaches—often extoled by software vendors—there is one that has widespread acceptance. A group of experts, software developers, consultants, and academics created the CRISP-DM Process in the late 1990s. Take a look at Figure 2-1 for a visual.



**Figure 2-1.** The CRISP-DM Process

In this chapter, we'll take a look at steps #1 through #3. Then in the rest of the book, we'll cover the remaining ones (that is, we will look at Modelling and Evaluation in Chapter 3 and Deployment in Chapter 8).

Note that steps #1–#3 can account for 80% of the time of the data process, which is based on the experience of Atif Kureishy, who is the global VP of Emerging Practices at Teradata.[14] This is due to factors like:

---

[14] This is from the author's interview with Atif Kureishy in February 2019.

The data is not well organized and comes from different sources (whether from different vendors or silos in the organization), there is not enough focus on automation tools, and the initial planning was insufficient for the scope of the project.

It's also worth keeping in mind that the CRISP-DM Process is not a strict linear process. When dealing with data, there can be much iteration. For example, there may be multiple attempts at coming up with the right data and testing it.

## Step #1—Business Understanding

You should come up with a clear view of the business problem to be solved. Some examples:

- How might a price adjustment impact your sales?

- Will a change in copy lead to improved conversion of digital ads?

- Does a fall in engagement mean there will be an increase in churn?

Then, you must establish how you will measure success. Might it be that sales should increase by at least 1% or that conversions should rise by 5%?

Here's a case from Prasad Vuyyuru, who is a partner at the Enterprise Insights Practice of Infosys Consulting:

> *Identifying which business problem to solve using AI and assessing what value will be created are critical for the success of all AI projects. Without such diligent focus on business value, AI projects risk not getting adopted in the organization. AB Inbev's experience in using AI to identify packaging line motors that are likely to fail is a great example of how AI is creating practical value. ABInbev installed 20 wireless sensors to measure vibrations at packaging lines motors. They compared sounds with normally functioning motors to identify anomalies which predicted eventual failure of the motors.[15]*

Regardless of the goal, it's essential that the process be free of any prejudgments or bias. The focus is to find the best results. No doubt, in some cases, there will not be a satisfactory result.

Or, in other situations, there may be big surprises. A famous example of this comes from the book *Moneyball* by Michael Lewis, which was also made into a movie in 2011 that starred Brad Pitt. It's a true story of how the Oakland A's used data science techniques to recruit players. The tradition in baseball

---

[15] This is from the author's interview of Prasad Vuyyuru in February 2019.

was to rely on metrics like batting averages. But when using sophisticated data analytics techniques, there were some startling results. The Oakland A's realized that the focus should be on slugging and on-base percentages. With this information, the team was able to recruit high-performing players at lower compensation levels.

The upshot is that you need to be open minded and willing to experiment.

In step #1, you should also assemble the right team for the project. Now unless you work at a company like Facebook or Google, you will not have the luxury of selecting a group of PhDs in machine learning and data science. Such talent is quite rare—and expensive.

But you also do not need an army of top-notch engineers for an AI project either. It is actually getting easier to apply machine learning and deep learning models, because of open source systems like TensorFlow and cloud-based platforms from Google, Amazon.com, and Microsoft. In other words, you may only need a couple people with a background in data science.

Next, you should find people—likely from your organization—who have the right domain expertise for the AI project. They will need to think through the workflows, models, and the training data—with a particular understanding of the industry and customer requirements.

Finally, you will need to evaluate the technical needs. What infrastructure and software tools will be used? Will there be a need to increase capacity or purchase new solutions?

# Step #2—Data Understanding

In this step, you will look at the data sources for the project. Consider that there are three main ones, which include the following:

- *In-House Data*: This data may come from a web site, beacons in a store location, IoT sensors, mobile apps, and so on. A major advantage of this data is that it is free and customized to your business. But then again, there are some risks. There can be problems if there has not been enough attention on the data formatting or what data should be selected.

- *Open Source Data*: This is usually freely available, which is certainly a nice benefit. Some examples of open source data include government and scientific information. The data is often accessed through an API, which makes the process fairly straightforward. Open source data is also usually well formatted. However, some of the variables may not be clear, and there could be bias, such as being skewed to a certain demographic.

- *Third-Party Data*: This is data from a commercial vendor. But the fees can be high. In fact, the data quality, in some cases, may be lacking.

According to Teradata—based on the firm's own AI engagements—about 70% of data sources are in-house, 20% from open source, and the rest from commercial vendors.[16] But despite the source, all data must be trusted. If not, there will likely be the problem of "garbage in, garbage out."

To evaluate the data, you need to answer questions like the following:

- Is the data complete? What might be missing?

- Where did the data come from?

- What were the collection points?

- Who touched the data and processed it?

- What have been the changes in the data?

- What are the quality issues?

If you are working with structured data, then this stage should be easier. However, when it comes to unstructured and semi-structured data, you will need to label the data—which can be a protracted process. But there are some tools emerging in the market that can help automate this process.

## Step #3—Data Preparation

The first step in the data preparation process is to decide what datasets to use.

Let's take a look at a scenario: Suppose you work for a publishing company and want to put together a strategy to improve customer retention. Some of the data that should help would include demographic information on the customer base like age, sex, income, and education. To provide more color, you can also look at browser information. What type of content interests customers? What's the frequency and duration? Any other interesting patterns—say accessing information during weekends? By combining the sources of information, you can put together a powerful model. For example, if there is a drop-off in activity in certain areas, it could pose a risk of cancellation. This would alert sales people to reach out to the customers.

---

[16] This is from the author's interview with Atif Kureishy in February 2019.

While this is a smart process, there are still landmines. Including or excluding even one variable can have a significant negative impact on an AI model. To see why, look back at the financial crisis. The models for underwriting mortgages were sophisticated and based on huge amounts of data. During normal economic times, they worked quite well as major financial institutions like Goldman Sachs, JP Morgan, and AIG relied on them heavily.

But there was a problem: The models did not account for falling housing prices! The main reason was that—for decades—there had never been a national drop. The assumption was that housing was mostly a local phenomenon.

Of course, housing prices more than just fell—they plunged. The models then proved to be far off the mark, and billions of dollars in losses nearly took down the US financial system. The federal government had little choice but to lend $700 billion for a bailout of Wall Street.

Granted, this is an extreme case. But it does highlight the importance of data selection. This is where having a solid team of domain experts and data scientists can be essential.

Next, when in the data preparation stage, there will need to be data cleansing. The fact is that all data has issues. Even companies like Facebook have gaps, ambiguities, and outliers in their datasets. It's inevitable.

So here are some actions you can take to cleanse the data:

- *De-duplication*: Set tests to identify any duplications and delete the extraneous data.

- *Outliers*: This is data that is well beyond the range of most of the rest of the data. This may indicate that the information is not helpful. But of course, there are situations where the reverse is true. This would be for fraud deduction.

- *Consistency*: Make sure you have clear definitions for the variables. Even terms like "revenue" or "customer" can have multiple meanings.

- *Validation Rules*: As you look at the data, try to find the inherent limitations. For example, you can have a flag for the age column. If it is over 120 in many cases, then the data has some serious issues.

- *Binning*: Certain data may not need to be specific. Does it really matter if someone is 35 or 37? Probably not. But comparing those from 30–40 to 41–50 probably would.

- *Staleness*: Is the data timely and relevant?

- *Merging*: In some cases, the columns of data may have very similar information. Perhaps one has height in inches and another in feet. If your model does not require a more detailed number, you can just use the one for feet.

- *One-Hot Encoding*: This is a way to replace categorical data as numbers. Example: Let's say we have a database with a column that has three possible values: Apple, Pineapple, and Orange. You could represent Apple as 1, Pineapple as 2, and Orange as 3. Sounds reasonable, right? Perhaps not. The problem is that an AI algorithm may think that Orange is greater than Apple. But with one-hot encoding, you can avoid this problem. You will create three new columns: is_Apple, is_Pineapple, and is_Orange. For each row in the data, you'll put 1 for where the fruit exists and 0 for the rest.

- *Conversion Tables*: You can use this when translating data from one standard to another. This would be the case if you have data in the decimal system and want to move over to the metric system.

These steps will go a long way in improving the quality of the data. There are also automation tools that can help out, such as from companies like SAS, Oracle, IBM, Lavastorm Analytics, and Talend. Then there are open source projects, such as OpenRefine, plyr, and reshape2.

Regardless, the data will not be perfect. No data source is. There will likely still be gaps and inaccuracies.

This is why you need to be creative. Look at what Eyal Lifshitz did, who is the CEO of BlueVine. His company leverages AI to provide financing to small businesses. "One of our data sources is credit information of our customers," he said. "But we've found that small business owners incorrectly identify their type of business. This could mean bad results for our underwriting. To deal with this, we scrape data from the customer website with AI algorithms, which helps identify the industry."[17]

Data cleansing approaches will also depend on the use cases for the AI project. For example, if you are building a system for predictive maintenance in manufacturing, the challenge will be to handle the wide variations from different sensors. The result is that a large amount of the data may have little value and be mostly noise.

---

[17] This is from the author's interview with Eyal Lifshitz in February 2019.

# Ethics and Governance

You need to be mindful of any restrictions on the data. Might the vendor prohibit you from using the information for certain purposes? Perhaps your company will be on the hook if something goes wrong? To deal with these issues, it is advisable to have the legal department brought in.

For the most part, data must be treated with care. After all, there are many high-profile cases where companies have violated privacy. A prominent example of this is Facebook. One of the company's partners, Cambridge Analytica, accessed millions of data points from profiles without the permission of users. When a whistleblower uncovered this, Facebook stock plunged—losing more than $100 billion in value. The company also came under pressure from the US and European governments.[18]

Something else to be wary of is scraping data from public sources. True, this is often an efficient way to create large datasets. There are also many tools that can automate the process. But scraping could expose your company to legal liability as the data may be subject to copyrights or privacy laws.

There are also some precautions that may ironically have inherent flaws. For example, a recent study from MIT shows that anonymized data may not be very anonymized. The researchers found that it was actually quite easy to reconstruct this type of data and identify the individuals—such as by merging two datasets. This was done by using data in Singapore from a mobile network (GPS tracking) and a local transportation system. After about 11 weeks of analysis, the researchers were able to identify 95% of the individuals.[19]

Finally, make sure you take steps to secure the data. The instances of cyberattacks and threats continue to increase at an alarming rate. In 2018, there were 53,000+ incidents and about 2,200 breaches, according to Verizon.[20] The report also noted the following:

- 76% of the breaches were financially motivated.

- 73% were from those outside the company.

- About half came from organized criminal groups and 12% from nation-state or state-affiliated actors.

The increasing use of cloud and on-premise data can subject a company to gaps in security as well. Then there is the mobile workforce, which can mean access to data that could expose it to breaches.

---

[18] https://venturebeat.com/2018/07/02/u-s-agencies-widen-investigation-into-what-facebook-knew-about-cambridge-analytica/
[19] http://news.mit.edu/2018/privacy-risks-mobility-data-1207
[20] https://enterprise.verizon.com/resources/reports/dbir/

The attacks are also getting much more damaging. The result is that a company can easily suffer penalties, lawsuits, and reputational damage.

Basically, when putting together an AI project, make sure there is a security plan and that it is followed.

# How Much Data Do You Need for AI?

The more data, the better, right? This is usually the case. Look at something called Hughes Phenomenon. This posits that as you add features to a model, the performance generally increases.

But quantity is not the end-all, be-all. There may come a point where the data starts to degrade. Keep in mind that you may run into something called the curse of dimensionality. According to Charles Isbell, who is the professor and senior associate dean of the School of Interactive Computing at Georgia Tech, "As the number of features or dimensions grows, the amount of data we need to generalize accurately grows exponentially."[21]

What is the practical impact? It could make it impossible to have a good model since there may not be enough data. This is why that when it comes to applications like vision recognition, the curse of dimensionality can be quite problematic. Even when analyzing RGB images, the number of dimensions is roughly 7,500. Just imagine how intensive the process would be using real-time, high-definition video.

# More Data Terms and Concepts

When engaging in data analysis, you should know the basic terms. Here are some that you'll often hear:

*Categorical Data*: This is data that does not have a numerical meaning. Rather, it has a textual meaning like a description of a group (race and gender). Although, you can assign numbers to each of the elements.

*Data Type*: This is the kind of information a variable represents, such as a Boolean, integer, string, or floating point number.

*Descriptive Analytics*: This is analyzing data to get a better understanding of the current status of a business. Some examples of this include measuring what products are selling better or determining risks in customer support. There are many traditional software tools for descriptive analytics, such as BI applications.

---

[21] www.kdnuggets.com/2017/04/must-know-curse-dimensionality.html

*Diagnostic Analytics*: This is querying data to see why something has happened. This type of analytics uses techniques like data mining, decision trees, and correlations.

*ETL (Extraction, Transformation, and Load)*: This is a form of data integration and is typically used in a data warehouse.

*Feature*: This is a column of data.

*Instance*: This is a row of data.

*Metadata*: This is data about data—that is, descriptions. For example, a music file can have metadata like the size, length, date of upload, comments, genre, artist, and so on. This type of data can wind up being quite useful for an AI project.

*Numerical Data*: This is any data that can be represented by a number. But numerical data can have two forms. There is discrete data, which is an integer—that is, a number without a decimal point. Then there is continuous data that has a flow, say temperature or time.

*OLAP (Online Analytical Processing)*: This is technology that allows you to analyze information from various databases.

*Ordinal Data*: This is a mix of numerical and categorical data. A common example of this is the five-star rating on Amazon.com. It has both a star and a number associated with it.

*Predictive Analytics*: This involves using data to make forecasts. The models for this are usually sophisticated and rely on AI approaches like machine learning. To be effective, it is important to update the underlying model with new data. Some of the tools for predictive analytics include machine learning approaches like regressions.

*Prescriptive Analytics*: This is about leveraging Big Data to make better decisions. This is not only focued on predicting outcomes—but understanding the rationales. And this is where AI plays a big part.

*Scalar Variables*: These are variables that hold single values like name or credit card number.

*Transactional Data*: This is data that is recorded on financial, business, and logistical actions. Examples include payments, invoices, and insurance claims.

# Conclusion

Being successful with AI means having a data-driven culture. This is what has been critical for companies like Amazon.com, Google, and Facebook. When making decisions, they look to the data first. There should also be wide availability of data across the organization.

Without this approach, success with AI will be fleeting, regardless of your planning. Perhaps this helps explain that—according to a study from NewVantage Partners—about 77% of respondents say that "business adoption" of Big Data and AI remain challenges.[22]

## Key Takeaways

- Structured data is labeled and formatted—and is often stored in a relational database or spreadsheet.

- Unstructured data is information that has no predefined formatting.

- Semi-structured data has some internal tags that help with categorization.

- Big Data describes a way to handle huge amounts of volumes of information.

- A relational database is based on relationships of data. But this structure can prove difficult for modern-day applications, such as AI.

- A NoSQL database is more free-form, being based on a document model. This has made it better able to deal with unstructured and semi-structured data.

- The CRISP-DM Process provides a way to manage data for a project, with steps that include business understanding, data understanding, data preparation, modelling, evaluation, and deployment.

- Quantity of data is certainly important, but there also needs to be much work on the quality. Even small errors can have a huge impact on the results of an AI model.

---

[22] http://newvantage.com/wp-content/uploads/2018/12/Big-Data-Executive-Survey-2019-Findings-Updated-010219-1.pdf

# Machine Learning

## Mining Insights from Data

*A breakthrough in machine learning would be worth ten Microsofts.*

—Bill Gates[1]

While Katrina Lake liked to shop online, she knew the experience could be much better. The main problem: It was tough to find fashions that were personalized.

So began the inspiration for Stitch Fix, which Katrina launched in her Cambridge apartment while attending Harvard Business School in 2011 (by the way, the original name for the company was the less catchy "Rack Habit"). The site had a Q&A for its users—asking about size and fashion styles, just to name a few factors—and expert stylists would then put together curated boxes of clothing and accessories that were sent out monthly.

The concept caught on quickly, and the growth was robust. But it was tough to raise capital as many venture capitalists did not see the potential in the business. Yet Katrina persisted and was able to create a profitable operation—fairly quickly.

---

[1] Steve Lohr, "Microsoft, Amid Dwindling Interest, Talks Up Computing as a Career: Enrollment in Computing Is Dwindling," *New York Times*, March 1, 2004, start page C1, quote page C2, column 6.

Along the way, Stitch Fix was collecting enormous amounts of valuable data, such as on body sizes and style preferences. Katrina realized that this would be ideal for machine learning. To leverage on this, she hired Eric Colson, who was the vice president of Data Science and Engineering at Netflix, his new title being chief algorithms officer.

This change in strategy was pivotal. The machine learning models got better and better with their predictions, as Stitch Fix collected more data—not only from the initial surveys but also from ongoing feedback. The data was also encoded in the SKUs.

The result: Stitch Fix saw ongoing improvement in customer loyalty and conversion rates. There were also improvements in inventory turnover, which helped to reduce costs.

But the new strategy did not mean firing the stylists. Rather, the machine learning greatly augmented their productivity and effectiveness.

The data also provided insights on what types of clothing to create. This led to the launch of Hybrid Designs in 2017, which is Stitch Fix's private-label brand. This proved effective in dealing with the gaps in inventory.

By November 2017, Katrina took Stitch Fix public, raising $120 million. The valuation of the company was a cool $1.63 billion—making her one of the richest women in the United States.[2] Oh, and at the time, she had a 14-month-old son!

Fast forward to today, Stitch Fix has 2.7 million customers in the United States and generates over $1.2 billion in revenues. There are also more than 100 data scientists on staff and a majority of them have PhDs in areas like neuroscience, mathematics, statistics, and AI.[3]

According to the company's 10-K filing:

> Our data science capabilities fuel our business. These capabilities consist of our rich and growing set of detailed client and merchandise data and our proprietary algorithms. We use data science throughout our business, including to style our clients, predict purchase behavior, forecast demand, optimize inventory and design new apparel.[4]

---

[2] www.cnbc.com/2017/11/16/stitch-fix-ipo-sees-orders-coming-in-under-range.html

[3] https://investors.stitchfix.com/static-files/2b398694-f553-4586-b763-e942617e4dbf

[4] www.sec.gov/Archives/edgar/data/1576942/000157694218000003/stitchfix201810k.htm

No doubt, the story of Stitch Fix clearly shows the incredible power of machine learning and how it can disrupt an industry. In an interview with digiday.com, Lake noted:

> Historically, there's been a gap between what you give to companies and how much the experience is improved. Big data is tracking you all over the web, and the most benefit you get from that right now is: If you clicked on a pair of shoes, you'll see that pair of shoes again a week from now. We'll see that gap begin to close. Expectations are very different around personalization, but importantly, an authentic version of it. Not, 'You abandoned your cart and we're recognizing that.' It will be genuinely recognizing who you are as a unique human. The only way to do this scalably is through embracing data science and what you can do through innovation.[5]

OK then, what is machine learning really about? Why can it be so impactful? And what are some of the risks to consider?

In this chapter, we'll answer these questions—and more.

## What Is Machine Learning?

After stints at MIT and Bell Telephone Laboratories, Arthur L. Samuel joined IBM in 1949 at the Poughkeepsie Laboratory. His efforts helped boost the computing power of the company's machines, such as with the development of the 701 (this was IBM's first commercialized computer system).

But he also programmed applications. And there was one that would make history—that is, his computer checkers game. It was the first example of a machine learning system (Samuel published an influential paper on this in 1959[6]). IBM CEO Thomas J. Watson, Sr., said that the innovation would add 15 points to the stock price![7]

Then why was Samuel's paper so consequential? By looking at checkers, he showed how machine learning works—in other words, a computer could learn and improve by processing data without having to be explicitly programmed. This was possible by leveraging advanced concepts of statistics, especially with probability analysis. Thus, a computer could be trained to make accurate predictions.

---

[5] https://digiday.com/marketing/stitch-fix-ceo-katrina-lake-predicts-ais-impact-fashion/
[6] Arthur L. Samuel, "Some Studies in Machine Learning Using the Game of Checkers," in Edward A. Feigenbaum and Julian Feldman, eds., *Computers and Thought* (New York: McGraw-Hill, 1983), pp. 71–105.
[7] https://history.computer.org/pioneers/samuel.html

This was revolutionary as software development, at this time, was mostly about a list of commands that followed a workflow of logic.

To get a sense of how machine learning works, let's use an example from the HBO TV comedy show *Silicon Valley*. Engineer Jian-Yang was supposed to create a Shazam for food. To train the app, he had to provide a massive dataset of food pictures. Unfortunately, because of time constraints, the app only learned how to identify…hot dogs. In other words, if you used the app, it would only respond with "hot dog" and "not hot dog."

While humorous, the episode did a pretty good job of demonstrating machine learning. In essence, it is a process of taking in labeled data and finding relationships. If you train the system with hot dogs—such as thousands of images—it will get better and better at recognizing them.

Yes, even TV shows can teach valuable lessons about AI!

But of course, you still need much more. In the next section of the chapter, we'll take a deeper look at the core statistics you need to know about machine learning. This includes the standard deviation, the normal distribution, Bayes' theorem, correlation, and feature extraction.

Then we'll cover topics like the use cases for machine learning, the general process, and the common algorithms.

# Standard Deviation

The standard deviation measures the average distance from the mean. In fact, there is no need to learn how to calculate this (the process involves multiple steps) since Excel or other software can do this for you easily.

To understand the standard deviation, let's take an example of the home values in your neighborhood. Suppose that the average is $145,000 and the standard deviation is $24,000. This means that one standard deviation below the average would be $133,000 ($145,000 − $12,000) and one standard deviation above the mean would come to $157,000 ($145,000 + $12,000). This gives us a way to quantify the variation in the data. That is, there is a spread of $24,000 from the average.

Next, let's take a look at the data if, well, Mark Zuckerberg moves into your neighborhood and, as a result, the average jumps to $850,000 and the standard deviation is $175,000. But do these statistical metrics reflect the valuations? Not really. Zuckerberg's purchase is an outlier. In this situation, the best approach may be instead to exclude his home.

# The Normal Distribution

When plotted on a graph, the normal distribution looks like a bell (this is why another name for it is the "bell curve"). It represents the sum of probabilities for a variable. Interestingly enough, the normal curve is common in the natural world, as it reflects distributions of such things like height and weight.

A general approach when interpreting a normal distribution is to use the 68-95-99.7 rule. This estimates that 68% of the data items will fall within one standard deviation, 95% within two standard deviations, and 99.7% within three standard deviations.

A way to understand this is to use IQ scores. Suppose the mean score is 100 and the standard deviation is 15. We'd have this for the three standard deviations, as shown in Figure 3-1.



**Figure 3-1.** Normal distribution of IQ scores

Note that the peak in this graph is the average. So, if a person has an IQ of 145, then only 0.15% will have a higher score.

Now the curve may have different shapes, depending on the variation in the data. For example, if our IQ data has a large number of geniuses, then the distribution will skew to the right.

# Bayes' Theorem

As the name implies, descriptive statistics provides information about your data. We've already seen this with such things as averages and standard deviations.

But of course, you can go well beyond this—basically, by using the Bayes' theorem. This approach is common in analyzing medical diseases, in which cause and effect are key—say for FDA (Federal Drug Administration) trials.

To understand how Bayes' theorem works, let's take an example. A researcher comes up with a test for a certain type of cancer, and it has proven to be accurate 80% of the time. This is known as a true positive.

But 9.6% of the time, the test will identify the person as having the cancer even though he or she does not have it, which is known as a false positive. Keep in mind that—in some drug tests—this percentage may be higher than the accuracy rate!

And finally, 1% of the population has the cancer.

In light of all this, if a doctor uses the test on you and it shows that you have the cancer, what is the probability that you really have the cancer? Well, Bayes' theorem will show the way. This calculation uses factors like accuracy rates, false positives, and the population rate to come up with a probability:

- *Step #1*: 80% accuracy rate × the chance of having the cancer (1%) = 0.008.

- *Step #2*: The chance of not having the cancer (99%) × the 9.6% false positive = 0.09504.

- *Step #3*: Then plug the above numbers into the following equation: 0.008 / (0.008 + 0.09504) = 7.8%.

Sounds kind of out of whack, right? Definitely. After all, how is it that a test, which is 90% accurate, has only a 7.8% probability of being right? But remember the accuracy rate is based on the measure of those who have the flu. And this is a small number since only 1% of the population has the flu. What's more, the test is still giving off false positives. So Bayes' theorem is a way to provide a better understanding of results—which is critical for systems like AI.

# Correlation

A machine learning algorithm often involves some type of correlation among the data. A quantitative way to describe this is to use the Pearson correlation, which shows the strength of the relationship between two variables that range from 1 to −1 (this is the coefficient).

Here's how it works:

- *Greater than 0*: This is where an increase in one variable leads to the increase in another. For example: Suppose that there is a 0.9 correlation between income and spending. If income increases by $1,000, then spending will be up by $900 ($1,000 × 0.9).

- *0*: There is no correlation between the two variables.

- *Less than 0*: Any increase in the variable means a decrease in another and vice versa. This describes an inverse relationship.

Then what is a strong correlation? As a general rule of thumb, it's if the coefficient is +0.7 or so. And if it is under 0.3, then the correlation is tenuous.

All this harkens the old saying of "Correlation is not necessarily causation." Yet when it comes to machine learning, this concept can easily be ignored and lead to misleading results.

For example, there are many correlations that are just random. In fact, some can be downright comical. Check out the following from Tylervigen.com:[8]

- The divorce rate in Maine has a 99.26% correlation with per capita consumption of margarine.

- The age of Miss America has an 87.01% correlation with the murders by steam, hot vapors, and hot tropics.

- The US crude oil imports from Norway have a 95.4% correlation with drivers killed in collision with a railway train.

There is a name for this: patternicity. This is the tendency to find patterns in meaningless noise.

# Feature Extraction

In Chapter 2, we looked at selecting the variables for a model. The process is often called feature extraction or feature engineering.

An example of this would be a computer model that identifies a male or female from a photo. For humans, this is fairly easy and quick. It's something that is intuitive. But if someone asked you to describe the differences, would you be able to? For most people, it would be a difficult task. However, if we want to build an effective machine learning model, we need to get feature extraction right—and this can be subjective.

---

[8] www.tylervigen.com/spurious-correlations

Table 3-1 shows some ideas about how a man's face may differ from a woman's.

**Table 3-1.** Facial features

| Features | Male |
| --- | --- |
| Eyebrows | Thicker and straighter |
| Face shape | Longer and larger, with more of a square shape |
| Jawbone | Square, wider, and sharper |
| Neck | Adam's apple |

This just scratches the surface as I'm sure you have your own ideas or approaches. And this is normal. But this is also why such things as facial recognition are highly complex and subject to error.

Feature extraction also has some nuanced issues. One is the potential for bias. For example, do you have preconceptions of what a man or woman looks like? If so, this can result in models that give wrong results.

Because of all this, it's a good idea to have a group of experts who can determine the right features. And if the feature engineering proves too complex, then machine learning is probably not a good option.

But there is another approach to consider: deep learning. This involves sophisticated models that find features in a data. Actually, this is one of the reasons that deep learning has been a major breakthrough in AI. We'll learn more about this in the next chapter.

# What Can You Do with Machine Learning?

As machine learning has been around for decades, there have been many uses for this powerful technology. It also helps that there are clear benefits, in terms of cost savings, revenue opportunities, and risk monitoring.

To give a sense of the myriad applications, here's a look at some examples:

- *Predictive Maintenance*: This monitors sensors to forecast when equipment may fail. This not only helps to reduce costs but also lessens downtime and boosts safety. In fact, companies like PrecisionHawk are actually using drones to collect data, which is much more efficient. The technology has proven quite effective for industries like energy, agriculture, and construction. Here's what PrecisionHawk notes about its own drone-based predictive maintenance system: "One client tested the use of visual line of sight (VLOS) drones to inspect a

cluster of 10 well pads in a three-mile radius. Our client determined that the use of drones reduced inspection costs by approximately 66%, from $80–$90 per well pad from traditional inspection methodology to $45–$60 per well pad using VLOS drone missions."[9]

- *Recruiting Employees*: This can be a tedious process since many resumes are often varied. This means it is easy to pass over great candidates. But machine learning can help in a big way. Take a look at CareerBuilder, which has collected and analyzed more than 2.3 million jobs, 680 million unique profiles, 310 million unique resumes, 10 million job titles, 1.3 billion skills, and 2.5 million background checks to build Hello to Hire. It's a platform that has leveraged machine learning to reduce the number of job applications—for a successful hire—to an average of 75. The industry average, on the other hand, is about 150.[10] The system also automates the creation of job descriptions, which even takes into account nuances based on the industry and location!

- *Customer Experience*: Nowadays, customers want a personalized experience. They have become accustomed to this by using services like Amazon.com and Uber. With machine learning, a company can leverage its data to gain insight—learning about what really works. This is so important that it led Kroger to buy a company in the space, called 84.51°. It is definitely key that it has data on more than 60 million US households. Here's a quick case study: For most of its stores, Kroger had bulk avocados, and only a few carried 4-packs. The conventional wisdom was that 4-packs had to be discounted because of the size disparity with the bulk items. But when applying machine learning analysis, this proved to be incorrect, as the 4-packs attracted new and different households like Millennials and ClickList shoppers. By expanding 4-packs across the chain, there was an overall increase in avocado sales.[11]

- *Finance*: Machine learning can detect discrepancies, say with billing. But there is a new category of technology, called RPA (Robotic Process Automation), that can help

---

[9] www.precisionhawk.com/blog/in-oil-gas-the-economics-of-bvlos-drone-operations

[10] This information is from the author's interview in February 2019 with Humair Ghauri, who is the chief product officer at CareerBuilder.

[11] www.8451.com/case-study/avocado

with this (we'll cover this topic in Chapter 5). It automates routine processes in order to help reduce errors. RPA also may use machine learning to detect abnormal or suspicious transactions.

- *Customer Service*: The past few years has seen the growth in chatbots, which use machine learning to automate interactions with customers. We'll cover this in Chapter 6.

- *Dating*: Machine learning could help find your soul mate! Tinder, one of the largest dating apps, is using the technology to help improve the matches. For instance, it has a system that automatically labels more than 10 billion photos that are uploaded on a daily basis.

Figure 3-2 shows some of the applications for machine learning.



**Figure 3-2.** Applications for machine learning

# The Machine Learning Process

To be successful with applying machine learning to a problem, it's important to take a systematic approach. If not, the results could be way off base.

First of all, you need to go through a data process, which we covered in the prior chapter. When this is finished, it's a good idea to do a visualization of the data. Is it mostly scattered? Or are there some patterns? If the answer is yes, then the data could be a good candidate for machine learning.

The goal of the machine learning process is to create a model, which is based on one or more algorithms. We develop this by training it. The goal is that the model should provide a high-degree of predictability.

Now let's take a closer look at this (by the way, this will also be applicable for deep learning, which we'll cover in the next chapter):

# Step #1—Data Order

If your data is sorted, then this could skew the results. That is, the machine learning algorithm may detect this as a pattern! This is why it's a good idea to randomize the order of the data.

# Step #2—Choose a Model

You will need to select an algorithm. This will be an educated guess, which will involve a process of trial and error. In this chapter, we'll look at the various algorithms available.

# Step #3—Train the Model

The training data, which will be about 70% of the complete dataset, will be used to create the relationships in the algorithm. For example, suppose you are building a machine learning system to find the value of a used car. Some of the features will include the year manufactured, make, model, mileage, and condition. By processing this training data, the algorithm will calculate the weights for each of these factors.

Example: Suppose we are using a linear regression algorithm, which has the following format:

$y = m * x + b$

In the training phase, the system will come up with the values for m (which is the slope on a graph) and b (which is the y-intercept).

# Step #4—Evaluate the Model

You will put together test data, which is the remaining 30% of the dataset. It should be representative of the ranges and type of information in the training data.

With the test data, you can see if the algorithm is accurate. In our used car example, are the market values consistent with what's happening in the real world?

---

■ **Note**    With the training and test data, there must not be any intermingling. This can easily lead to distorted results. Interestingly enough, this is a common mistake.

---

Now accuracy is one measure of the success of the algorithm. But this can, in some cases, be misleading. Consider the situation with fraud deduction. There are usually a small number of features when compared to a dataset. But missing one could be devastating, costing a company millions of dollars in losses.

This is why you might want to use other approaches like Bayes' theorem.

## Step #5—Fine-Tune the Model

In this step, we can adjust the values of the parameters in the algorithm. This is to see if we can get better results.

When fine-tuning the model, there may also be hyperparameters. These are parameters that cannot be learned directly from the training process.

# Applying Algorithms

Some algorithms are quite easy to calculate, while others require complex steps and mathematics. The good news is that you usually do not have to compute an algorithm because there are a variety of languages like Python and R that make the process straightforward.

As for machine learning, an algorithm is typically different from a traditional one. The reason is that the first step is to process data—and then, the computer will start to learn.

Even though there are hundreds of machine learning algorithms available, they can actually be divided into four major categories: supervised learning, unsupervised learning, reinforcement learning, and semi-supervised learning. We'll take a look at each.

## Supervised Learning

Supervised learning uses labeled data. For example, suppose we have a set of photos of thousands of dogs. The data is considered to be labeled if each photo identifies each for the breed. For the most part, this makes it easier to analyze since we can compare our results with the correct answer.

One of the keys with supervised learning is that there should be large amounts of data. This helps to refine the model and produce more accurate results.

But there is a big issue: The reality is that much of the data available is not labeled. In addition, it could be time consuming to provide labels if there is a massive dataset.

Yet there are creative ways to deal with this, such as with crowdfunding. This is how the ImageNet system was built, which was a breakthrough in AI innovation. But it still took several years to create it.

Or, in some cases, there can be automated approaches to label data. Take the example of Facebook. In 2018, the company announced—at its F8 developers conference—it leveraged its enormous database of photos from Instagram, which were labeled with hashtags.[12]

Granted, this approach had its flaws. A hashtag may give a nonvisual description of the photo—say #tbt (which is "throwback Thursday)—or could be too vague, like #party. This is why Facebook called its approach "weakly supervised data." But the talented engineers at the company found some ways to improve the quality, such as by building a sophisticated hashtag prediction model.

All in all, things worked out quite well. Facebook's machine learning model, which included 3.5 billion photos, had an accuracy rate of 85.4%, which was based on the ImageNet recognition benchmark. It was actually the highest recorded in history, by 2%.

This AI project also required innovative approaches for building the infrastructure. According to the Facebook blog post:

> Since a single machine would have taken more than a year to complete the model training, we created a way to distribute the task across up to 336 GPUs, shortening the total training time to just a few weeks. With ever-larger model sizes—the biggest in this research is a ResNeXt 101-32x48d with over 861 million parameters—such distributed training is increasingly essential. In addition, we designed a method for removing duplicates to ensure we don't accidentally train our models on images that we want to evaluate them on, a problem that plagues similar research in this area.[13]

Going forward, Facebook sees potential in using its approach to various areas, including the following:

- Improved ranking in the newsfeed

- Better detection of objectionable content

- Auto generation of captions for the visually impaired

---

[12] www.engadget.com/2018/05/02/facebook-trained-image-recognition-ai-instagram-pics/
[13] https://code.fb.com/ml-applications/advancing-state-of-the-art-image-recognition-with-deep-learning-on-hashtags/

# Unsupervised Learning

Unsupervised learning is when you are working with unlabeled data. This means that you will use deep learning algorithms to detect patterns.

By far, the most common approach for unsupervised learning is clustering, which takes unlabeled data and uses algorithms to put similar items into groups. The process usually starts with guesses, and then there are iterations of the calculations to get better results. At the heart of this is finding data items that are close together, which can be accomplished with a variety of quantitative methods:

- *Euclidean Metric*: This is a straight line between two data points. The Euclidean metric is quite common with machine learning.

- *Cosine Similarity Metric*: As the name implies, you will use a cosine to measure the angle. The idea is to find similarities between two data points in terms of the orientation.

- *Manhattan Metric*: This involves taking the sum of the absolute distances of two points on the coordinates of a graph. It's called the "Manhattan" because it references the city's street layout, which allows for shorter distances for travel.

In terms of use cases for clustering, one of the most common is customer segmentation, which is to help better target marketing messages. For the most part, a group that has similar characteristics is likely to share interests and preferences.

Another application is sentiment analysis, which is where you mine social media data and find the trends. For a fashion company, this can be crucial in understanding how to adapt the styles of the upcoming line of clothes.

Now there are other approaches than just clustering. Here's a look at three more:

- *Association*: The basic concept is that if X happens, then Y is likely to happen. Thus, if you buy my book on AI, you will probably want to buy other titles in the genre. With association, a deep learning algorithm can decipher these kinds of relationships. This can result in powerful recommendation engines.

- *Anomaly Detection*: This identifies outliers or anomalous patterns in the dataset, which can be helpful with cybersecurity applications. According to Asaf Cidon, who is the VP of Email Security at Barracuda Networks: "We've found that by combining many different signals—such as the email body, header, the social graph of communications, IP logins, inbox forwarding rules, etc.—we're able to achieve an extremely high precision in detecting social engineering attacks, even though the attacks are highly personalized and crafted to target a particular person within a particular organization. Machine learning enables us to detect attacks that originate from within the organization, whose source is a legitimate mailbox of an employee, which would be impossible to do with a static one-size-fits-all rule engine."[14]

- *Autoencoders*: With this, the data will be put into a compressed form, and then it will be reconstructed. From this, new patterns may emerge. However, the use of autoencoders is rare. But it could be shown to be useful in helping with applications like reducing noise in data.

Consider that many AI researchers believe that unsupervised learning will likely be critical for the next level of achievements. According to a paper in *Nature* by Yann LeCun, Geoffrey Hinton, and Yoshua Bengio, "We expect unsupervised learning to become far more important in the longer term. Human and animal learning is largely unsupervised: we discover the structure of the world by observing it, not by being told the name of every object."[15]

## Reinforcement Learning

When you were a kid and wanted to play a new sport, chances were you did not read a manual. Instead, you observed what other people were doing and tried to figure things out. In some situations, you made mistakes and lost the ball as your teammates would show their displeasure. But in other cases, you made the right moves and scored. Through this trial-and-error process, your learning was improved based on positive and negative reinforcement.

---

[14] This is from the author's interview in February 2019 with Asaf Cidon, who is the VP of Email Security at Barracuda Networks.
[15] https://towardsdatascience.com/simple-explanation-of-semi-supervised-learning-and-pseudo-labeling-c2218e8c769b

At a high level, this is analogous to reinforcement learning. It has been key for some of the most notable achievements in AI, such as the following:

- *Games*: They are ideal for reinforcement learning since there are clear-cut rules, scores, and various constraints (like a game board). When building a model, you can test it with millions of simulations, which means that the system will quickly get smarter and smarter. This is how a program can learn to beat the world champion of Go or chess.

- *Robotics*: A key is being able to navigate within a space— and this requires evaluating the environment at many different points. If the robot wants to move to, say, the kitchen, it will need to navigate around furniture and other obstacles. If it runs into things, there will be a negative reinforcement action.

## Semi-supervised Learning

This is a mix of supervised and unsupervised learning. This arises when you have a small amount of unlabeled data. But you can use deep learning systems to translate the unsupervised data to supervised data—a process that is called pseudo-labeling. After this, you can then apply the algorithms.

An interesting use case of semi-supervised learning is the interpretation of MRIs. A radiologist can first label the scans, and after this, a deep learning system can find the rest of the patterns.

# Common Types of Machine Learning Algorithms

There is simply not enough room in this book to cover all the machine learning algorithms! Instead, it's better to focus on the most common ones.

In the remaining part of this chapter, we'll take a look at those for the following:

- *Supervised Learning*: You can boil down the algorithms to two variations. One is classification, which divides the dataset into common labels. Examples of the algorithms include Naive Bayes Classifier and k-Nearest Neighbor (neural networks will be covered in Chapter 4). Next, there is regression, which finds continuous patterns in the data. For this, we'll take a look at linear regression, ensemble modelling, and decision trees.

- *Unsupervised Learning*: In this category, we'll look at clustering. For this, we'll cover k-Means clustering.

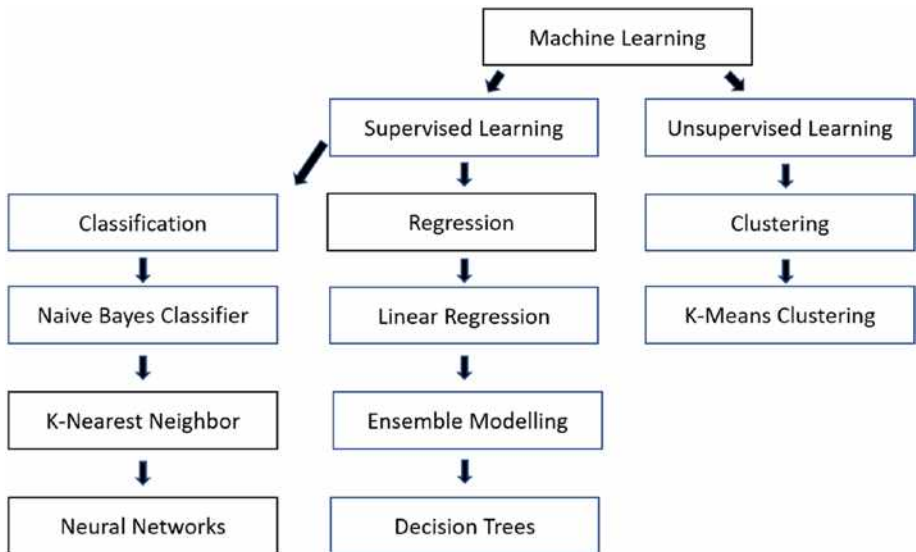Figure 3-3 shows a general framework for machine learning algorithms.



**Figure 3-3.** General framework for machine learning algorithms

# Naïve Bayes Classifier (Supervised Learning/Classification)

Earlier in this chapter, we looked at Bayes' theorem. As for machine learning, this has been modified into something called the Naïve Bayes Classifier. It is "naïve" because the assumption is that the variables are independent from each other—that is, the occurrence of one variable has nothing to do with the others. True, this may seem like a drawback. But the fact is that the Naïve Bayes Classifier has proven to be quite effective and fast to develop.

There is another assumption to note as well: the a priori assumption. This says that the predictions will be wrong if the data has changed.

There are three variations on the Naïve Bayes Classifier:

- *Bernoulli*: This is if you have binary data (true/false, yes/no).

- *Multinomial*: This is if the data is discrete, such as the number of pages of a book.

- *Gaussian*: This is if you are working with data that conforms to a normal distribution.

A common use case for Naïve Bayes Classifiers is text analysis. Examples include email spam detection, customer segmentation, sentiment analysis, medical diagnosis, and weather predictions. The reason is that this approach is useful in classifying data based on key features and patterns.

To see how this is done, let's take an example: Suppose you run an e-commerce site and have a large database of customer transactions. You want to see how variables like product review ratings, discounts, and time of year impact sales.

Table 3-2 shows a look at the dataset.

**Table 3-2.** Customer transactions dataset

| Discount | Product Review | Purchase |
|----------|----------------|----------|
| Yes | High | Yes |
| Yes | Low | Yes |
| No | Low | No |
| No | Low | No |
| No | Low | No |
| No | High | Yes |
| Yes | High | No |
| Yes | Low | Yes |
| No | High | Yes |
| Yes | High | Yes |
| No | High | No |
| No | Low | Yes |
| Yes | High | Yes |
| Yes | Low | No |

You will then organize this data into frequency tables, as shown in Tables 3-3 and 3-4.

**Table 3-3.** Discount frequency table

| | | Purchase | |
|----------|-----|-----|-----|
| | | **Yes** | **No** |
| **Discount** | Yes | 19 | 1 |
| | Yes | 5 | 5 |

**Table 3-4.** Product review frequency table

| | | Purchase | | |
|---|---|---|---|---|
| | | Yes | No | Total |
| **Product Review** | **High** | 21 | 2 | 11 |
| | **Low** | 3 | 4 | 8 |
| | **Total** | 24 | 6 | 19 |

When looking at this, we call the purchase an event and the discount and product reviews as independent variables. Then we can make a probability table for one of the independent variables, say the product reviews. See Table 3-5.

**Table 3-5.** Product review probability table

| | | Purchase | | |
|---|---|---|---|---|
| | | Yes | No | |
| **Product Reviews** | **High** | 9/24 | 2/6 | 11/30 |
| | **Low** | 7/24 | 1/6 | 8/30 |
| | | 24/30 | 6/30 | |

Using this chart, we can see that the probability of a purchase when there is a low product review is 7/24 or 29%. In other words, the Naïve Bayes Classifier allows more granular predictions within a dataset. It is also relatively easy to train and can work well with small datasets.

# K-Nearest Neighbor (Supervised Learning/Classification)

The k-Nearest Neighbor (k-NN) is a method for classifying a dataset (k represents the number of neighbors). The theory is that those values that are close together are likely to be good predictors for a model. Think of it as "Birds of a feather flock together."

A use case for k-NN is the credit score, which is based on a variety of factors like income, payment histories, location, home ownership, and so on. The algorithm will divide the dataset into different segments of customers. Then, when there is a new customer added to the base, you will see what cluster he or she falls into—and this will be the credit score.

K-NN is actually simple to calculate. In fact, it is called lazy learning because there is no training process with the data.

To use k-NN, you need to come up with the distance between the nearest values. If the values are numerical, it could be based on a Euclidian distance, which involves complicated math. Or, if there is categorical data, then you can use an overlap metric (this is where the data is the same or very similar).

Next, you'll need to identify the number of neighbors. While having more will smooth the model, it can also mean a need for huge amount of computational resources. To manage this, you can assign higher weights to data that are closer to their neighbors.

# Linear Regression (Supervised Learning/Regression)

Linear regression shows the relationship between certain variables. The equation—assuming there is enough quality data—can help predict outcomes based on inputs.

Example: Suppose we have data on the number of hours spent studying for an exam and the grade. See Table 3-6.

**Table 3-6.** Chart for hours of study and grades

| Hours of Study | Grade Percentage |
| --- | --- |
| I | 0.75 |
| I | 0.69 |
| I | 0.71 |
| 3 | 0.82 |
| 3 | 0.83 |
| 4 | 0.86 |
| 5 | 0.85 |
| 5 | 0.89 |
| 5 | 0.84 |
| 6 | 0.91 |
| 6 | 0.92 |
| 7 | 0.95 |

As you can see, the general relationship is positive (this describes the tendency where a higher grade is correlated with more hours of study). With the regression algorithm, we can plot a line that has the best fit (this is done by using a calculation called "least squares," which minimizes the errors). See Figure 3-4.
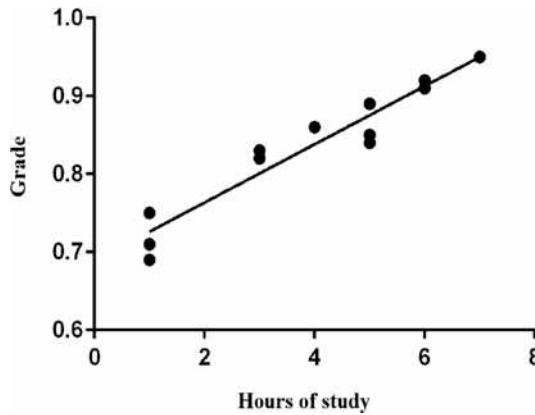
**Figure 3-4.** This is a plot of a linear regression model that is based on hours of study

From this, we get the following equation:

Grade = Number of hours of study × 0.03731 + 0.6889

Then, let's suppose you study 4 hours for the exam. What will be your estimated grade? The equation tells us how:

0.838 = 4 × 0.03731 + 0.6889

How accurate is this? To help answer this question, we can use a calculation called R-squared. In our case, it is 0.9180 (this ranges from 0 to 1). The closer the value is to 1, the better the fit. So 0.9180 is quite high. It means that the hours of study explains 91.8% of the grade on the exam.

Now it's true that this model is simplistic. To better reflect reality, you can add more variables to explain the grade on the exam—say the student's attendance. When doing this, you will use something called multivariate regression.

---

■ **Note**   If the coefficient for a variable is quite small, then it might be a good idea to not include it in the model.

---

Sometimes data may not be in a straight line either, in which case the regression algorithm will not work. But you can use a more complex version, called polynomial regression.

# Decision Tree (Supervised Learning/Regression)

No doubt, clustering may not work on some datasets. But the good news is that there are alternatives, such as a decision tree. This approach generally works better with nonnumerical data.

The start of a decision tree is the root node, which is at the top of the flow chart. From this point, there will be a tree of decision paths, which are called splits. At these points, you will use an algorithm to make a decision, and there will be a probability computed. At the end of the tree will be the leaf (or the outcome).

A famous example—in machine learning circles—is to use a decision tree for the tragic sinking of the Titanic. The model predicts the survival of a passenger based on three features: sex, age, and the number of spouses or children along (sibsp). Here's how it looks, in Figure 3-5.



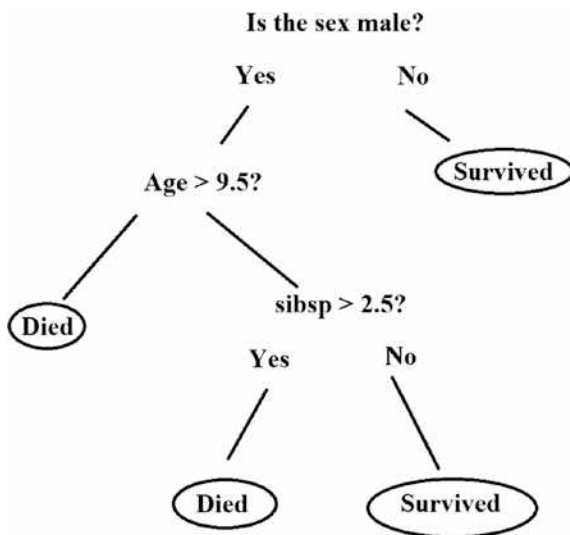**Figure 3-5.** This is a basic decision tree algorithm for predicting the survival of the Titanic

There are clear advantages for decision trees. They are easy to understand, work well with large datasets, and provide transparency with the model.

However, decision trees also have drawbacks. One is error propagation. If one of the splits turns out to be wrong, then this error can cascade throughout the rest of the model!

Next, as the decision trees grow, there will be more complexity as there will be a large number of algorithms. This could ultimately result in lower performance for the model.

# Ensemble Modelling (Supervised Learning/Regression)

Ensemble modelling means using more than one model for your predictions. Even though this increases the complexity, this approach has been shown to generate strong results.

To see this in action, take a look at the "Netflix Prize," which began in 2006. The company announced it would pay $1 million to anyone or any team that could improve the accuracy of its movie recommendation system by 10% or more. Netflix also provided a dataset of over 100 million ratings of 17,770 movies from 480,189 users.[16] There would ultimately be more than 30,000 downloads.

Why did Netflix do all this? A big reason is that the company's own engineers were having trouble making progress. Then why not give it to the crowd to figure out? It turned out to be quite ingenious—and the $1 million payout was really modest compared to the potential benefits.

The contest certainly stirred up a lot of activity from coders and data scientists, ranging from students to employees at companies like AT&T.

Netflix also made the contest simple. The main requirement was that the teams had to disclose their methods, which helped boost the results (there was even a dashboard with rankings of the teams).

But it was not until 2009 that a team—BellKor's Pragmatic Chaos—won the prize. Then again, there were considerable challenges.

So how did the winning team pull it off? The first step was to create a baseline model that smoothed out the tricky issues with the data. For example, some movies only had a handful of ratings, whereas others had thousands. Then there was the thorny problem where there were users who would always rate a movie with one star. To deal with these matters, BellKor used machine learning to predict ratings in order to fill the gaps.

---

[16] www.thrillist.com/entertainment/nation/the-netflix-prize

Once the baseline was finished, there were more tough challenges to tackle like the following:

- A system may wind up recommending the same films to many users.

- Some movies may not fit well within genres. For example, *Alien* is really a cross of science fiction and horror.

- There were movies, like *Napoleon Dynamite*, that proved extremely difficult for algorithms to understand.

- Ratings of a movie would often change over time.

The winning team used ensemble modelling, which involved hundreds of algorithms. They also used something called boosting, which is where you build consecutive models. With this, the weights in the algorithms are adjusted based on the results of the previous model, which help the predictions get better over time (another approach, called bagging, is when you build different models in parallel and then select the best one).

But in the end, BellKor found the solutions. However, despite this, Netflix did not use the model! Now it's not clear why this was the case. Perhaps it was that Netflix was moving away from five-star ratings anyway and was more focused on streaming. The contest also had blowback from people who thought there may have been privacy violations.

Regardless, the contest did highlight the power of machine learning—and the importance of collaboration.

# K-Means Clustering (Unsupervised/Clustering)

The k-Means clustering algorithm, which is effective for large datasets, puts similar, unlabeled data into different groups. The first step is to select k, which is the number of clusters. To help with this, you can perform visualizations of that data to see if there are noticeable grouping areas.

Here's a look at sample data, in Figure 3-6:



**Figure 3-6.** The initial plot for a dataset

For this example, we assume there will be two clusters, and this means there will also be two centroids. A centroid is the midpoint of a cluster. We will assign each randomly, which you can see in Figure 3-7.
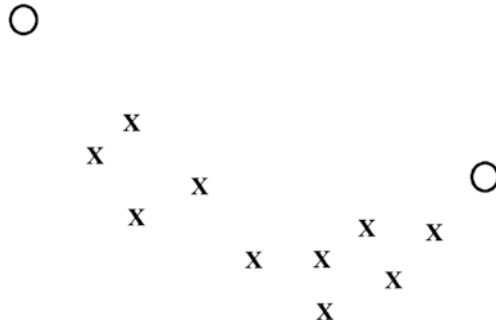


**Figure 3-7.** This chart shows two centroids—represented by circles—that are randomly placed

As you can see, the centroid at the top left looks way off, but the one on the right side is better. The k-Means algorithm will then calculate the average distances of the centroids and then change their locations. This will be iterated until the errors are fairly minimal—a point that is called convergence, which you can see with Figure 3-8.
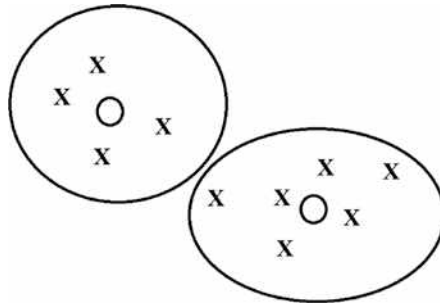


**Figure 3-8.** Through iterations, the k-Means algorithm gets better at grouping the data

Granted, this is a simple illustration. But of course, with a complex dataset, it will be difficult to come up with the number of initial clusters. In this situation, you can experiment with different k values and then measure the average distances. By doing this multiple times, there should be more accuracy.

Then why not just have a high number for k? You can certainly do this. But when you compute the average, you'll notice that there will be only incremental improvements. So one method is to stop at the point where this starts to occur. This is seen in Figure 3-9.

**Figure 3-9.** This shows the optimal point of the k value in the k-Means algorithm

However, k-Means has its drawbacks. For instance, it does not work well with nonspherical data, which is the case with Figure 3-10.



**Figure 3-10.** Here's a demonstration where k-Means does not work with nonspherical data

With this, the k-Means algorithm would likely not pick up on the surrounding data, even though it has a pattern. But there are some algorithms that can help, such as DBScan (density-based spatial clustering of applications with

noise), which is meant to handle a mix of widely varying sizes of datasets. Although, DBScan can require lots of computational power.

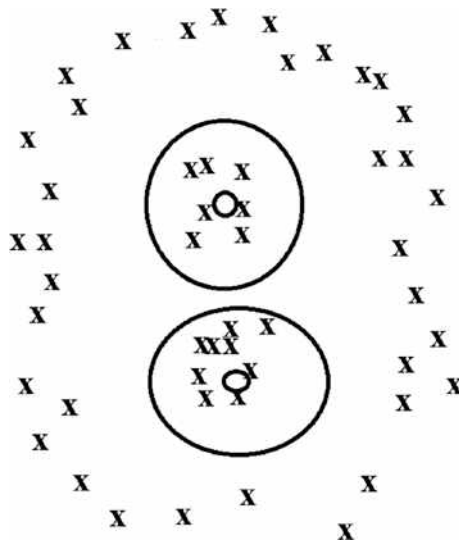Next, there is the situation where there are some clusters with lots of data and others with little. What might happen? There is a chance that the k-Means algorithm will not pick up on the light one. This is the case with Figure 3-11.



**Figure 3-11.** *If there are areas of thin data, the k-Means algorithm may not pick them up*

# Conclusion

These algorithms can get complicated and do require strong technical skills. But it is important to not get too bogged down in the technology. After all, the focus is to find ways to use machine learning to accomplish clear objectives.

Again, Stich Fix is a good place to get guidance on this. In the November issue of the *Harvard Business Review*, the company's chief algorithms officer, Eric Colson, published an article, "Curiosity-Driven Data Science."[17] In it, he provided his experiences in creating a data-driven organization.

At the heart of this is allowing data scientists to explore new ideas, concepts, and approaches. This has resulted in AI being implemented across core functions of the business like inventory management, relationship management, logistics, and merchandise buying. It has been transformative, making the

---

[17] https://hbr.org/2018/11/curiosity-driven-data-science

organization more agile and streamlined. Colson also believes it has provided "a protective barrier against competition."

His article also provides other helpful advice for data analysis:

- *Data Scientists*: They should not be part of another department. Rather, they should have their own, which reports directly to the CEO. This helps with focusing on key priorities as well as having a holistic view of the needs of the organization.

- *Experiments*: When a data scientist has a new idea, it should be tested on a small sample of customers. If there is traction, then it can be rolled out to the rest of the base.

- *Resources*: Data scientists need full access to data and tools. There should also be ongoing training.

- *Generalists*: Hire data scientists who span different domains like modelling, machine learning, and analytics (Colson refers to these people as "full-stack data scientists"). This leads to small teams—which are often more efficient and productive.

- *Culture*: Colson looks for values like "learning by doing, being comfortable with ambiguity, balancing long-and short-term returns."

## Key Takeaways

- Machine learning, whose roots go back to the 1950s, is where a computer can learn without being explicitly programmed. Rather, it will ingest and process data by using sophisticated statistical techniques.

- An outlier is data that is far outside the rest of the numbers in the dataset.

- The standard deviation measures the average distance from the mean.

- The normal distribution—which has a shape like a bell—represents the sum of probabilities for a variable.

- The Bayes' theorem is a sophisticated statistical technique that provides a deeper look at probabilities.

- A true positive is when a model makes a correct prediction. A false positive, on the other hand, is when a model prediction shows that the result is true even though it is not.

- The Pearson correlation shows the strength of the relationship between two variables that range from 1 to -1.

- Feature extraction or feature engineering describes the process of selecting variables for a model. This is critical since even one wrong variable can have a major impact on the results.

- Training data is what is used to create the relationships in an algorithm. The test data, on the other hand, is used to evaluate the model.

- Supervised learning uses labeled data to create a model, whereas unsupervised learning does not. There is also semi-supervised learning, which uses a mix of both approaches.

- Reinforcement learning is a way to train a model by rewarding accurate predictions and punishing those that are not.

- The k-Nearest Neighbor (k-NN) is an algorithm based on the notion that values that are close together are good predictors for a model.

- Linear regression estimates the relationship between certain variables. The R-squared will indicate the strength of the relationship.

- A decision tree is a model that is based on a workflow of yes/no decisions.

- An ensemble model uses more than one model for the predictions.

- The k-Means clustering algorithm puts similar unlabeled data into different groups.

# Deep Learning

## The Revolution in AI

*Take any old classification problem where you have a lot of data, and it's going to be solved by deep learning. There's going to be thousands of applications of deep learning.*

—Geoffrey Hinton,
English Canadian cognitive psychologist and computer scientist[1]

Fei-Fei Li, who got a BA degree in physics from Princeton in 1999 with high honors and a PhD in electrical engineering from Caltech in 2005, focused her brilliance on developing AI models. But she had a major challenge: finding quality datasets. At first, she looked at creating them by hand, such as with graduate students who downloaded images from the Internet. But the process was too slow and tedious.

One day a student mentioned to Li that Amazon.com's Mechanical Turk, an online service that uses crowdsourcing to solve problems, could be a good way to scale the process. It would allow for fast and accurate labeling of the data.

Li gave it a try, and it worked out quite well. By 2010, she had created ImageNet, which had 3.2 million images across over 5,200 categories.

---

[1] Siddhartha Mukherjee, "The Algorithm Will See You Now," *The New Yorker*, April 3, 2017, https://www.newyorker.com/magazine/2017/04/03/ai-versus-md.

Yet it got a tepid response from the academic community. But this did not deter Li. She continued to work tirelessly to evangelize the dataset. In 2012, she put together a contest as a way to encourage researchers to create more effective models and push the boundaries of innovation. It would turn out to be a game changer, and the contest would become an annual event.

In the first contest, professors from the University of Toronto—Geoffrey Hinton, Ilya Sutskever, and Alex Krizhevsky—used sophisticated deep learning algorithms. And the results were standout. The system they built, which was called AlexNet, beat all the other contestants by a margin of 10.8%.[2]

This was no fluke. In the years after this, deep learning continued to show accelerated progress with ImageNet. As of now, the error rate for deep learning is a mere 2% or so—which is better than humans.

By the way, Li has since gone on to become a professor at Stanford and co-director of the school's AI lab. She is also Google's chief scientist of AI and Machine Learning. Needless to say, whenever she has new ideas now, people listen!

In this chapter, we'll take a look at deep learning, which is clearly the hottest area of AI. It has led to major advances in areas like self-driving cars and virtual assistants like Siri.

Yes, deep learning can be a complicated subject, and the field is constantly changing. But we'll take a look at the main concepts and trends—without getting into the technical details.

# Difference Between Deep Learning and Machine Learning

There is often confusion between deep learning and machine learning. And this is reasonable. Both topics are quite complex, and they do share many similarities.

So to understand the differences, let's first take a look at two high-level aspects of machine learning and how they relate to deep learning. First of all, while both usually require large amounts of data, the types are generally different.

Take the following example: Suppose we have photos of thousands of animals and want to create an algorithm to find the horses. Well, machine learning cannot analyze the photos themselves; instead, the data must be labeled. The machine learning algorithm will then be trained to recognize horses, through a process known as supervised learning (covered in Chapter 3).

---

[2]https://qz.com/1034972/the-data-that-changed-the-direction-of-ai-research-and-possibly-the-world/

Even though machine learning will likely come up with good results, they will still have limitations. Wouldn't it be better to look at the pixels of the images themselves—and find the patterns? Definitely.

But to do this with machine learning, you need to use a process called feature extraction. This means you must come up with the kinds of characteristics of a horse—such as the shape, the hooves, color, and height—which the algorithms will then try to identify.

Again, this is a good approach—but it is far from perfect. What if your features are off the mark or do not account for outliers or exceptions? In such cases, the accuracy of the model will likely suffer. After all, there are many variations to a horse. Feature extraction also has the drawback of ignoring a large amount of the data. This can be exceedingly complicated—if not impossible—for certain use cases. Look at computer viruses. Their structures and patterns, which are known as signatures, are constantly changing so as to infiltrate systems. But with feature extraction, a person would somehow have to anticipate this, which is not practical. This is why cybersecurity software is often about collecting signatures after a virus has exacted damage.

But with deep learning, we can solve these problems. This approach analyzes all the data—pixel by pixel—and then finds the relationships by using a neural network, which mimics the human brain.

Let's take a look.

## So What Is Deep Learning Then?

Deep learning is a subfield of machine learning. This type of system allows for processing huge amounts of data to find relationships and patterns that humans are often unable to detect. The word "deep" refers to the number of hidden layers in the neural network, which provide much of the power to learn.

When it comes to the topic of AI, deep learning is at the cutting-edge and often generates most of the buzz in mainstream media. "[Deep learning] AI is the new electricity," extolled Andrew Yan-Tak Ng, who is the former chief scientist at Baidu and co-founder of Google Brain.[3]

But it is also important to remember that deep learning is still in the early stages of development and commercialization. For example, it was not until about 2015 that Google started using this technology for its search engine.

---

[3] https://medium.com/@GabriellaLeone/the-best-explanation-machine-learning-vs-deep-learning-d5c123405b11

As we saw in Chapter 1, the history of neural networks was full of ebbs and flows. It was Frank Rosenblatt who created the perceptron, which was a fairly basic system. But real academic progress with neural networks did not occur until the 1980s, such as with the breakthroughs with backpropagation, convolutional neural networks, and recurrent neural networks. But for deep learning to have an impact on the real world, it would take the staggering growth in data, such as from the Internet, and the surge in computing power.

# The Brain and Deep Learning

Weighing only about 3.3 pounds, the human brain is an amazing feat of evolution. There are about 86 billion neurons—often called gray matter—that are connected with trillions of synapses. Think of neurons as CPUs (Central Processing Units) that take in data. The learning occurs with the strengthening or weakening of the synapses.

The brain is made up of three regions: the forebrain, the midbrain, and the hindbrain. Among these, there are a variety of areas that perform different functions. Some of the main ones include the following:

- *Hippocampus*: This is where your brain stores memories. In fact, this is the part that fails when a person has Alzheimer's disease, in which a person loses the ability to form short-term memories.

- *Frontal Lobe*: Here the brain focuses on emotions, speech, creativity, judgment, planning, and reasoning.

- *Cerebral Cortex*: This is perhaps the most important when it comes to AI. The cerebral cortex helps with thinking and other cognitive activities. According to research from Suzana Herculano-Houzel, the level of intelligence is related to the number of neurons in this area of the brain.

Then how does deep learning compare to the human brain? There are some tenuous similarities. At least in areas like the retina, there is a process of ingesting data and processing them through a complex network, which is based on assigning weights. But of course, this is only a minute part of the learning process. Besides, there are still many mysteries about the human brain, and of course, it is not based on things like digital computing (instead, it appears that it is more of an analogue system). However, as the research continues to advance, the discoveries in neuroscience could help build new models for AI.

# Artificial Neural Networks (ANNs)

At the most basic level, an artificial neural network (ANN) is a function that includes units (which may also be called neurons, perceptrons, or nodes). Each unit will have a value and a weight, which indicates the relative importance, and will go into the hidden layer. The hidden layer uses a function, with the result becoming the output. There is also another value, called bias, which is a constant and is used in the calculation of the function.

This type of training of a model is called a feed-forward neural network. In other words, it only goes from input to the hidden layer to the output. It does not cycle back. But it could go to a new neural network, with the output becoming the input.

Figure 4-1 shows a chart of a feed-forward neural network.



**Figure 4-1.** A basic feed-forward neural network

Let's go deeper on this by taking an example. Suppose you are creating a model to predict whether a company's stock will increase. The following are what the variables represent as well as the values and weights assigned:

- $X_1$: Revenues are growing at a minimum of 20% a year. The value is 2.

- $X_2$: The profit margin is at least 20%. The value is 4.

- $W_1$: 1.9.

- $W_2$: 9.6.

- b: This is the bias (the value is 1), which helps smooth out the calculations.

You'll then sum the weights, and then the function will process the information. This will often involve an activation function, which is non-linear. This is more reflective of the real world since data is usually not in a straight line.

Now there are a variety of activation functions to choose from. One of the most common is the sigmoid. This compresses the input value into a range of 0–1. The closer it is to 1, the more accurate the model.

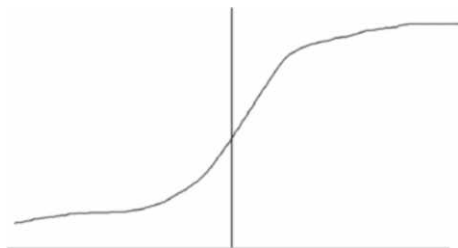When you graph this function, it will look like an S shape. See Figure 4-2.



**Figure 4-2.** A typical sigmoid activation function

As you can see, the system is relatively simplistic and will not be helpful in high-end AI models. To add much more power, there usually needs to be multiple hidden layers. This results in a multilayered perceptron (MLP). It also helps to use something called backpropagation, which allows for the output to be circled back into the neural network.

# Backpropagation

One of the major drawbacks with artificial neural networks is the process of making adjustments to the weights in the model. Traditional approaches, like the use of the mutation algorithm, used random values that proved to be time consuming.

Given this, researchers looked for alternatives, such as backpropagation. This technique had been around since the 1970s but got little interest as the performance was lacking. But David Rumelhart, Geoffrey Hinton, and Ronald Williams realized that backpropagation still had potential, so long as it was refined. In 1986, they wrote a paper entitled "Learning Representations by Back-propagating Errors," and it was a bombshell in the AI community.[4] It clearly showed that backpropagation could be much faster but also allow for more powerful artificial neural networks.

As should be no surprise, there is a lot of math involved in backpropagation. But when you boil things down, it's about adjusting the neural network when errors are found and then iterating the new values through the neural network again. Essentially, the process involves slight changes that continue to optimize the model.

---

[4] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams, "Learning Representations by Back-propagating Errors," *Nature* 323 (1986): 533–536.

For example, let's say one of the inputs has an output of 0.6. This means that the error is 0.4 (1.0 minus 0.6), which is subpar. But we can then backpropogate the output, and perhaps the new output may get to 0.65. This training will go on until the value is much closer to 1.

Figure 4-3 illustrates this process. At first, there is a high level of errors because the weights are too large. But by making iterations, the errors will gradually fall. However, doing too much of this could mean an increase in errors. In other words, the goal of backpropagation is to find the midpoint.



**Figure 4-3.** The optimal value for a backpropagation function is at the bottom of the graph

As a gauge of the success of backpropagation, there were a myriad of commercial applications that sprung up. One was called NETtalk, which was developed by Terrence Sejnowski and Charles Rosenberg in the mid-1980s. The machine was able to learn how to pronounce English text. NETtalk was so interesting that it was even demoed on the *Today show*.

There were also a variety of startups that were created that leveraged backpropagation, such as HNC Software. It built models that detected credit card fraud. Up until this point—when HNC was founded in the late 1980s—the process was done mostly by hand, which led to costly errors and low volumes of issuances. But by using deep learning approaches, credit card companies were able to save billions of dollars.

In 2002, HNC was acquired by Fair, Isaac and valued at $810 million.[5]

---

[5] www.insurancejournal.com/news/national/2002/05/01/16857.htm

# The Various Neural Networks

The most basic type of a neural network is a fully connected neural network. As the name implies, it is where all the neurons have connections from layer to layer. This network is actually quite popular since it means having to use little judgment when creating the model.

Then what are some of the other neural networks? The common ones include the recurrent neural network (RNN), the convolutional neural network (CNN), and the generative adversarial network (GAN), which we'll cover next.

## Recurrent Neural Network

With a recurrent neural network (RNN), the function not only processes the input but also prior inputs across time. An example of this is what happens when you enter characters in a messaging app. As you begin to type, the system will predict the words. So if you tap out "He," the computer will suggest "He," "Hello," and "Here's." The RNN is essentially a string of neural networks that feed on each other based on complex algorithms.

There are variations on the model. One is called LSTM, which stands for long short-term memory. This came about from a paper written by professors Sepp Hochreiter and Jürgen Schmidhuber in 1997.[6] In it, they set forth a way to effectively use inputs that are separated from each other for long time periods, allowing the use of more datasets.

Of course, RNNs do have drawbacks. There is the vanishing gradient problem, which means that the accuracy decays as the models get larger. The models can also take longer to train.

To deal with this, Google developed a new model called the Transformer, which is much more efficient since it processes the inputs in parallel. It also results in more accurate results.

Google has gained much insight about RNNs through its Translate app, which handles over 100 languages and processes over 100 billion words a day.[7] Launched in 2006, it initially used machine learning systems. But in 2016, Google switched to deep learning by creating Google Neural Machine Translation.[8] All in all, it has resulted in much higher accuracy rates.[9]

---

[6] Sepp Hochreiter and Jürgen Schmidhuber, "Long Short-Term Memory," *Neural Computation* 9, no. 8 (1997): 1735-80.
[7] www.argotrans.com/blog/accurate-google-translate-2018/
[8] www.techspot.com/news/75637-google-translate-not-monetized-despite-converting-over-100.html
[9] www.argotrans.com/blog/accurate-google-translate-2018/

Consider how Google Translate has helped out doctors who work with patients who speak other languages. According to a study from the University of California, San Francisco (UCSF), that was published in *JAMA Internal Medicine*, the app had a 92% accuracy rate with English-to-Spanish translations. This was up from 60% over the past couple years.[10]

# Convolutional Neural Network (CNN)

Intuitively, it makes sense to have all the units in a neural network to be connected. This works well with many applications.

But there are scenarios where it is far from optimal, such as with image recognition. Just imagine how complex a model would be where every pixel is a unit! It could quickly become unmanageable. There would also be other complications like overfitting. This is where the data is not reflective of what is being tested or there is a focus on the wrong features.

To deal with all this, you can use a convolutional neural network (CNN). The origins of this go back to professor Yann LeCun in 1998, when he published a paper called "Gradient-Based Learning Applied to Document Recognition."[11] Despite its strong insights and breakthroughs, it got little traction. But as deep learning started to show significant progress in 2012, researchers revisited the model.

LeCun got his inspiration for the CNN from Nobel Prize winners David Hubel and Torsten Wiesel who studied neurons of the visual cortex. This system takes an image from the retina and processes it in different stages—from easy to more complex. Each of the stages is called a convolution. For example, the first level would be to identify lines and angles; next, the visual cortex will find the shapes; and then it will detect the objects.

This is analogous to how a computer-based CNN works. Let's take an example: Suppose you want to build a model that can identify a letter. The CNN will have input in the form of an image that has 3,072 pixels. Each of the pixels will have a value that is from 0 to 255, which indicates the overall intensity. By using a CNN, the computer will go through multiple variations to identify the features.

The first is the convolutional layer, which is a filter that scans the image. In our example, this could be 5 × 5 pixels. The process will create a feature map, which is a long array of numbers. Next, the model will apply more filters to the image. By doing this, the CNN will identify the lines, edges and shapes—all

---

[10] https://gizmodo.com/google-translate-can-help-doctors-bridge-the-language-g-1832881294
[11] Yann LeCun et al., "Gradient-Based Learning Applied to Document Recognition," *Proceedings of the IEEE* 86 no. 11 (1998): 2278-2324.

expressed in numbers. With the various output layers, the model will use pooling, which combines them to generate a single output, and then create a fully connected neural network.

A CNN can definitely get complex. But it should be able to accurately identify the numbers that are input into the system.

## Generative Adversarial Networks (GANs)

Ian Goodfellow, who got his masters in computer science at Stanford and his PhD in machine learning at the Université de Montréal, would go on to work at Google. In his 20s, he co-authored one of the top books in AI, called *Deep Learning*,[12] and also made innovations with Google Maps.

But it was in 2014 that he had his most impactful breakthrough. It actually happened in a pub in Montreal when he talked with some of his friends about how deep learning could *create* photos.[13] At the time, the approach was to use generative models, but they were often blurry and nonsensical.

Goodfellow realized that there had to be a better why. So why not use game theory? That is, have two models compete against each other in a tight feedback loop. This could also be done with unlabeled data.

Here's a basic workflow:

- *Generator*: This neural network creates a myriad of new creations, such as photos or sentences.

- *Discriminator*: This neural network looks at the creations to see which ones are real.

- *Adjustments*: With the two results, a new model would change the creations to make them as realistic as possible. Through many iterations, the discriminator will no longer need to be used.

He was so excited about the idea that after he left the pub he started to code his ideas. The result was a new deep learning model: the generative adversarial network or GAN. And the results were standout. He would soon become an AI rock star.

---

[12] Ian Goodfellow, Yoshua Bengio, and Aaron Courville, *Deep Learning* (Cambridge, MA: The MIT Press, 2016).
[13] www.technologyreview.com/s/610253/the-ganfather-the-man-whos-given-machines-the-gift-of-imagination/

GAN research has already spurred over 500 academic papers.[14] Companies like Facebook have also used this technology, such as for its photo analysis and processing. The company's chief AI scientist, Yann LeCun, noted that GANs are the "the coolest idea in deep learning in the last 20 years."[15]

GANs have also been shown to help with sophisticated scientific research. For example, they have helped improve the accuracy of detecting behavior of subatomic particles in the Large Hadron Collider at CERN in Switzerland.[16]

While still in the early innings, this technology could lead to such things as a computer that can develop new types of fashion items or maybe a new-fangled wearable. Perhaps a GAN could even come up with a hit rap song.

And it could be sooner than you think. As a teenager, Robbie Barrat taught himself how to use deep learning systems and built a model to rap in the style of Kanye West.

But this was just the beginning of his AI wizardry. As a researcher at Stanford, he developed his own GAN platform, which processed roughly 10,000 nude portraits. The system then would create truly mesmerizing new works of art (you can find them at his Twitter account at @DrBeef_).

Oh, and he also made his system open source at his GitHub account. This caught the attention of a collective of French artists, called Obvious, that used the technology to create portraits of an eighteenth-century fictional family. It was based on processing 15,000 portraits from the fourteenth to the twentieth centuries.

In 2018, Obvious put its artwork at a Christie's auction, fetching a cool $432,000. [17]

But unfortunately, when it comes to GANs, there have been uses that have been less than admirable. One example is to use them for deepfakes, which involve leveraging neural networks to create images or videos that are misleading. Some of this is just kind of playful. For example, one GAN makes it possible to have Barack Obama say anything you tell him!

Yet there are lots of risks. Researchers at New York University and the Michigan State University wrote a paper that focused on "DeepMasterPrints."[18]

---

[14] https://github.com/hindupuravinash/the-gan-zoo
[15] https://trendsandevents4developers.wordpress.com/2017/04/24/the-coolest-idea-in-deep-learning-in-20-years-and-more/
[16] www.hpcwire.com/2018/08/14/cern-incorporates-ai-into-physics-based-simulations/
[17] www.washingtonpost.com/nation/2018/10/26/year-old-developed-code-ai-portrait-that-sold-christies/?utm_term=.b2f366a4460e
[18] www.cnbc.com/2018/12/28/research-claims-fake-fingerprints-could-hack-a-third-of-smartphones.html

It showed how a GAN can develop fake fingerprints to unlock three types of smartphones!

Then there was the incident of a so-called deepfake video of actress Jennifer Lawrence at a Golden Globes press conference. Her face was merged with Steve Buscemi's.[19]

# Deep Learning Applications

With so much money and resources being devoted to deep learning, there has been a surge in innovations. It seems that every day there is something amazing that is being announced.

Then what are some of the applications? Where has deep learning proven to be a game changer? Let's take a look at some that cover areas like healthcare, energy, and even earthquakes

## Use Case: Detecting Alzheimer's Disease

Despite decades of research, a cure for Alzheimer's disease remains elusive. Although, scientists have developed drugs that have slowed down the progression of the disease.

In light of this, early diagnosis is critical—and deep learning can potentially be a big help. Researchers at the UCSF Department of Radiology and Biomedical Imaging have used this technology to analyze brain screens—from the Alzheimer's Disease Neuroimaging Initiative public dataset—and to detect changes in the levels of glucose.

The result: The model can diagnose Alzheimer's disease up to six years before a clinical diagnosis. One of the tests showed a 92% accuracy rate, and another was 98%.

Now this is still in the beginning phases—and there will need to be more datasets analyzed. But so far, the results are very encouraging.

According to Dr. Jae Ho Sohn, who authored the study:

> This is an ideal application of deep learning because it is particularly strong at finding very subtle but diffuse processes. Human radiologists are really strong at identifying tiny focal finding like a brain tumor, but we struggle at detecting more slow, global changes. Given the strength of deep learning in this type of application, especially compared to humans, it seemed like a natural application.[20]

---

[19] http://fortune.com/2019/01/31/what-is-deep-fake-video/
[20] www.ucsf.edu/news/2018/12/412946/artificial-intelligence-can-detect-alzheimers-disease-brain-scans-six-years

# Use Case: Energy

Because of its massive data center infrastructure, Google is one of the largest consumers of energy. Even a small improvement in efficiency can lead to a sizeable impact on the bottom line. But there could also be the benefits of less carbon emissions.

To help with these goals, Google's DeepMind unit has been applying deep learning, which has involved better management of wind power. Even though this is a clean source of energy, it can be tough to use because of the changes in weather.

But DeepMind's deep learning algorithms have been critical. Applied to 700 megawatts of wind power in the United States, they were able to make accurate forecasts for output with a lead time of 36 hours.

According to DeepMind's blog:

> This is important, because energy sources that can be scheduled (i.e. can deliver a set amount of electricity at a set time) are often more valuable to the grid…To date, machine learning has boosted the value of our wind energy by roughly 20 percent, compared to the baseline scenario of no time-based commitments to the grid.[21]

But of course, this deep learning system could be more than just about Google—it could have a wide-ranging impact on energy use across the world.

# Use Case: Earthquakes

Earthquakes are extremely complicated to understand. They are also exceedingly difficult to predict. You need to evaluate faults, rock formations and deformations, electromagnetic activity, and changes in the groundwater. Hey, there is even evidence that animals have the ability to sense an earthquake!

But over the decades, scientists have collected huge amounts of data on this topic. In other words, this could be an application for deep learning, right?

Absolutely.

Seismologists at Caltech, which include Yisong Yue, Egill Hauksson, Zachary Ross, and Men-Andrin Meier, have been doing considerable research on this, using convolutional neural networks and recurrent neural networks. They are trying to build an effective early-warning system.

---

[21] https://deepmind.com/blog/machine-learning-can-boost-value-wind-energy/

Here's what Yue had to say:

> AI can [analyze earthquakes] faster and more accurately than humans can, and even find patterns that would otherwise escape the human eye. Furthermore, the patterns we hope to extract are hard for rule-based systems to adequately capture, and so the advanced pattern-matching abilities of modern deep learning can offer superior performance than existing automated earthquake monitoring algorithms.[22]

But the key is improving data collection. This means more analysis of small earthquakes (in California, there is an average of 50 each day). The goal is to create an earthquake catalog that can lead to the creation of a virtual seismologist, who can make evaluations of an earthquake faster than a human. This could allow for faster lead times when an earthquake strikes, which may help to save lives and property.

## Use Case: Radiology

PET scans and MRIs are amazing technology. But there are definitely downsides. A patient needs to stay within a confining tube for 30 minutes to an hour. This is uncomfortable and means being exposed to gadolinium, which has been shown to have harmful side effects.

Greg Zaharchuk and Enhao Gong, who met at Stanford, thought there could be a better way. Zaharchuk was an MD and PhD, with a specialization in radiology. He was also the doctoral advisor of Gong, who was an electrical engineering PhD in deep learning and medical image reconstruction.

In 2017, they co-founded Subtle Medical and hired some of the brightest imaging scientists, radiologists, and AI experts. Together, they set themselves to the challenge of improving PET scans and MRIs. Subtle Medical created a system that not only reduces the time for an MRI and PET scans by up to ten times, but the accuracy has been much higher. This was powered by high-end NVIDIA GPUs.

Then in December 2018, the system received FDA (Federal Drug Administration) 510(k) clearance and a CE mark approval for the European market.[23] It was the first ever AI-based nuclear medical device to achieve both of these designations.

---

[22] www.caltech.edu/about/news/qa-creating-virtual-seismologist-84789
[23] https://subtlemedical.com/subtle-medical-receives-fda-510k-clearance-and-ce-mark-approval-for-subtlepet/

Subtle Medical has more plans to revolutionize the radiology business. As of 2019, it is developing SubtleMRTM, which will be even more powerful than the company's current solution, and SubtleGADTM, which will reduce gadolinium dosages.[24]

# Deep Learning Hardware

Regarding chip systems for deep learning, GPUs have been the primary choice. But as AI gets more sophisticated—such as with GANs—and the datasets much larger, there is certainly more room for new approaches. Companies also have custom needs, such as in terms of functions and data. After all, an app for a consumer is usually quite different than one that is focused on the enterprise.

As a result, some of the mega tech companies have been developing their own chipsets:

- *Google*: In the summer of 2018, the company announced its third version of its Tensor Processing Unit (TPU; the first chip was developed in 2016).[25] The chips are so powerful—handling over 100 petaflops for training of models—there needs to be liquid cooling in the data centers. Google has also announced a version of its TPU for devices. Essentially, it means that processing will have less latency because there will be no need to access the cloud.

- *Amazon*: In 2018, the company announced AWS Inferentia.[26] The technology, which has come out of the acquisition of Annapurna in 2015, is focused on handling complex inference operations. In other words, this is what happens after a model has been trained.

- *Facebook and Intel*: These companies have joined forces to create an AI chip.[27] But the initiative is still in the initial phases. Intel also has been getting traction with an AI chip called the Nervana Neural Network Processor (NNP).

---

[24] www.streetinsider.com/Press+Releases/Subtle+Medical+Receives+FDA+510%28k%29+Clearance+and+CE+Mark+Approval+for+SubtlePET™/14892974.html
[25] www.theregister.co.uk/2018/05/09/google_tpu_3/
[26] https://aws.amazon.com/about-aws/whats-new/2018/11/announcing-amazon-inferentia-machine-learning-inference-microchip/
[27] www.analyticsindiamag.com/inference-chips-are-the-next-big-battlefield-for-nvidia-and-intel/

- *Alibaba*: The company has created its own AI chip company called Pingtouge.[28] It also has plans to build a quantum computer processor, which is based on qubits (they represent subatomic particles like electrons and photons).

- *Tesla*: Elon Musk has developed his own AI chip. It has 6 billion transistors and can process 36 trillion operations per second.[29]

There are a variety of startups that are making a play for the AI chip market as well. Among the leading companies is Untether AI, which is focused on creating chips that boost the transfer speeds of data (this has been a particularly difficult part of AI). In one of the company's prototypes, this process was more than 1,000 faster than a typical AI chip.[30] Intel, along with other investors, participated in a $13 million round of funding in 2019.

Now when it comes to AI chips, NVIDIA has the dominant market share. But because of the importance of this technology, it seems inevitable that there will be more and more offerings that will come to the market.

## When to Use Deep Learning?

Because of the power of deep learning, there is the temptation to use this technology first when creating an AI project. But this can be a big mistake. Deep learning still has narrow use cases, such as for text, video, image, and time-series datasets. There is also a need for large amounts of data and high-powered computer systems.

Oh, and deep learning is better when outcomes can be quantified and verified.

To see why, let's consider the following example. A team of researchers, led by Thomas Hartung (a toxicologist at Johns Hopkins University), created a dataset of about 10,000 chemicals that were based on 800,000 animal tests. By using deep learning, the results showed that the model was more predictive than many animal tests for toxicity.[31] Remember that animal tests can not only be costly and require safety measures but also have inconsistent results because of repeated testing on the same chemical.

---

[28] www.technologyreview.com/s/612190/why-alibaba-is-investing-in-ai-chips-and-quantum-computing/

[29] www.technologyreview.com/f/613403/tesla-says-its-new-self-driving-chip-will-help-make-its-cars-autonomous/

[30] www.technologyreview.com/f/613258/intel-buys-into-an-ai-chip-that-can-transfer-data-1000-times-faster/

[31] www.nature.com/articles/d41586-018-05664-2

"The first scenario illustrates the predictive power of deep learning, and its ability to unearth correlations from large datasets that a human would never find," said Sheldon Fernandez, who is the CEO of DarwinAI.[32]

So where's a scenario in which deep learning falls short? Actually, an illustration of this is the 2018 FIFA World Cup in Russia, which France won. Many researchers tried to predict the outcomes of all 64 matches, but the results were far from accurate:[33]

- One group of researchers employed the bookmaker consensus model that indicated that Brazil would win.

- Another group of researchers used algorithms such as random forest and Poisson ranking to forecast that Spain would prevail.

The problem here is that it is tough to find the right variables that have predictive power. In fact, deep learning models are basically unable to handle the complexity of features for certain events, especially those that have elements of being chaotic.

However, even if you have the right amount of data and computing power, you still need to hire people who have a background in deep learning, which is not easy. Keep in mind that it is a challenge to select the right model and fine-tune it. How many hyperparameters should there be? What should be the number of hidden layers? And how do you evaluate the model? All of these are highly complex.

Even experts can get things wrong. Here's the following from Sheldon:

One of our automotive clients encountered some bizarre behavior in which a self-driving car would turn left with increasing regularity when the sky was a certain shade of purple. After months of painful debugging, they determined the training for certain turning scenarios had been conducted in the Nevada desert when the sky was a particular hue. Unbeknownst to its human designers, the neural network had established a correlation between its turning behavior and the celestial tint.[34]

There are some tools that are helping with the deep learning process, such as Amazon.com's SageMaker, Google's HyperTune, and SigOpt. But there is still a long way to go.

If deep learning is not a fit, then you may want to consider machine learning, which often requires relatively less data. Furthermore, the models tend to be much simpler, but the results may still be more effective.

---

[32] This is from the author's interview with Sheldon Fernandez, the CEO of DarwinAI.
[33] https://medium.com/futuristone/artificial-intelligence-failed-in-world-cup-2018-6af10602206a
[34] This is from the author's interview with Sheldon Fernandez, the CEO of DarwinAI.

# Drawbacks with Deep Learning

Given all the innovations and breakthroughs, it's reasonable that many people consider deep learning to be a silver bullet. It will mean we no longer have to drive a car. It may even mean that we'll cure cancer.

How is it not possible to be excited and optimistic? This is natural and reasonable. But it is important to note that deep learning is still in a nascent stage and there are actually many nagging issues. It's a good idea to temper expectations.

In 2018, Gary Marcus wrote a paper entitled "Deep Learning: A Critical Appraisal," in which he clearly set forth the challenges.[35] In his paper, he notes:

> Against a background of considerable progress in areas such as speech recognition, image recognition, and game playing, and considerable enthusiasm in the popular press, I present ten concerns for deep learning, and suggest that deep learning must be supplemented by other techniques if we are to reach Artificial General Intelligence.[36]

Marcus definitely has the right pedigree to present his concerns, as he has both an academic and business background in AI. Before becoming a professor at the Department of Psychology at New York University, he sold his startup, called Geometric Intelligence, to Uber. Marcus is also the author of several bestselling books like *The Haphazard Construction of the Human Mind*.[37]

Here's a look at some of his worries about deep learning:

- *Black Box*: A deep learning model could easily have millions of parameters that involve many hidden layers. Having a clear understanding of this is really beyond a person's capabilities. True, this may not necessarily be a problem with recognizing cats in a dataset. But it could definitely be an issue with models for medical diagnosis or determining the safety of an oil rig. In these situations, regulators will want to have a good understanding of the transparency of the models. Because of this, researchers are looking at creating systems to determine "explainability," which provides an understanding of deep learning models.

---

[35] Gary Marcus, "Deep Learning: A Critical Appraisal," arXiv, 1801.00631v1 [cs.AI]:1–27, 2018.

[36] https://arxiv.org/ftp/arxiv/papers/1801/1801.00631.pdf

[37] Gary Marcus, *Kluge: The Haphazard Construction of the Human Mind* (Houghton Mifflin, 2008).

- *Data*: The human brain has its flaws. But there are some functions that it does extremely well like the ability to learn by abstraction. For example, suppose Jan, who is five years old, goes to a restaurant with her family. Her mother points out an item on the plate and says it is a "taco." She does not have to explain it or provide any information about it. Instead, Jan's brain will instantly process this information and understand the overall pattern. In the future, when she sees another taco—even if it has differences, such as with the dressing—she will know what it is. For the most part, this is intuitive. But unfortunately, when it comes to deep learning, there is no taco learning by abstraction! The system has to process enormous amounts of information to recognize it. Of course, this is not a problem for companies like Facebook, Google, or even Uber. But many companies have much more limited datasets. The result is that deep learning may not be a good option.

- *Hierarchical Structure*: This way of organizing does not exist in deep learning. Because of this, language understanding still has a long way to go (especially with long discussions).

- *Open-Ended Inference*: Marcus notes that deep learning cannot understand the nuances between "John promised Mary to leave" and "John promised to leave Mary." What's more, deep learning is far away from being able to, for instance, read Jane Austen's *Pride and Prejudice* and be able to divine Elizabeth Bennet's character motivations.

- *Conceptual Thinking*: Deep learning cannot have an understanding of concepts like democracy, justice, or happiness. It also does not have imagination, thinking of new ideas or plans.

- *Common Sense*: This is something deep learning does not do well. If anything, this means a model can be easily confused. For example, let's say you ask an AI system, "Is it possible to make a computer with a sponge?" For the most part, it will probably not know that this is a ridiculous question.

- *Causation*: Deep learning is unable to determine this. It's all about finding correlations.

- *Prior Knowledge*: CNNs can help with some prior information, but this is limited. Deep learning is still fairly self-contained, as it only solves one problem at a time. It cannot take in the data and create algorithms that span various domains. In addition, a model does not adapt. If there is change in the data, then a new model needs to be trained and tested. And finally, deep learning does not have prior understanding of what people know instinctively—such as basic physics of the real world. This is something that has to be explicitly programmed into an AI system.

- *Static*: Deep learning works best in environments that are fairly simple. This is why AI has been so effective with board games, which have a clear set of rules and boundaries. But the real world is chaotic and unpredictable. This means that deep learning may fall short with complex problems, even with self-driving cars.

- *Resources*: A deep learning model often requires a tremendous amount of CPU power, such as with GPUs. This can get costly. Although, one option is to use a third-party cloud service.

This is quite a lot? It's true. But the paper still has left out some drawbacks. Here are a couple other ones:

- *Butterfly Effect*: Because of the complexity of the data, networks, and connections, a minute change can have a major impact in the results of the deep learning model. This could easily lead to conclusions that are wrong or misleading.

- *Overfitting*: We explained this concept earlier in the chapter.

As for Marcus, his biggest fear is that AI could "get trapped in a local minimum, dwelling too heavily in the wrong part of intellectual space, focusing too much on the detailed exploration of a particular class of accessible but limited models that are geared around capturing low-hanging fruit—potentially neglecting riskier excursions that might ultimately lead to a more robust path."

However, he is not a pessimist. He believes that researchers need to go beyond deep learning and find new techniques that can solve tough problems.

# Conclusion

While Marcus has pointed out the flaws in deep learning, the fact is that this AI approach is still extremely powerful. In less than a decade, it has revolutionized the tech world—and is also significantly impacting areas like finance, robotics, and healthcare.

With the surge in investments from large tech companies and VCs, there will be further innovation with the models. This will also encourage engineers to get postgraduate degrees, creating a virtuous cycle of breakthroughs.

# Key Takeaways

- Deep learning, which is a subfield of machine learning, processes huge amounts of data to detect relationships and patterns that humans are often unable to detect. The word "deep" describes the number of hidden layers.

- An artificial neural network (ANN) is a function that includes units that have weights and are used to predict values in an AI model.

- A hidden layer is a part of a model that processes incoming data.

- A feed-forward neural network has data that goes only from input to the hidden layer to the output. The results do not cycle back. Yet they can go into another neural network.

- An activation function is non-linear. In other words, it tends to do a better job of reflecting the real world.

- A sigmoid is an activation function that compresses the input value into a range of 0–1, which makes it easier for analysis.

- Backpropagation is a sophisticated technique to adjust the weights in a neural network. This approach has been critical for the growth in deep learning.

- A recurrent neural network (RNN) is a function that not only processes the input but also prior inputs across time.

- A convolutional neural network (CNN) analyzes data section by section (that is, by convolutions). This model is geared for complex applications like image recognition.

- A generative adversarial network or GAN is where two neural networks compete with each other in a tight feedback loop. The result is often the creation of a new object.

- Explainability describes techniques for transparency with complex deep learning models.

# Robotic Process Automation (RPA)

## An Easier Path to AI

*By interacting with applications just as a human would, software robots can open email attachments, complete e-forms, record and re-key data, and perform other tasks that mimic human action.*

—Kaushik Iyengar,
director of Digital Transformation and Optimization at AT&T[1]

Back in 2005, Daniel Dines and Marius Tirca founded UiPath, which was located in Bucharest, Romania. The company focused mostly on providing integration services for applications from Google, Microsoft, and IBM. But it was a struggle as the company relied mostly on custom work for clients.

---

[1] www2.deloitte.com/insights/us/en/focus/signals-for-strategists/cognitive-enterprise-robotic-process-automation.html

By 2013, UiPath was close to being shut down. But the founders did not give up as they saw this as an opportunity to rethink the business and find a new opportunity.[2] To this end, they started to build a platform for Robotic Process Automation (RPA). The category, which had been around since 2000, was about automating routine and mundane tasks within a company.

Yet RPA was actually a backwater area in the tech world—as seen with the slow growth rates. However, Dines and Tirca were convinced that they could transform the industry. One of the key reasons: the rise of AI and the cloud.

The new strategy was spot-on, and growth took off. Dines and Tirca also were aggressive with seeking funding, innovating its RPA platform, and expanding into global markets.

By 2018, UiPath was considered the fastest-growing enterprise software company—ever. The annual recurring revenue soared from $1 million to $100 million, with over 1,800 customers.[3] The company had the most widely used RPA system in the world.

UiPath attracted a total of $448 million in venture capital from marque firms like CapitalG, Sequoia Capital, and Accel. The valuation was at $3 billion.

In light of all this, more RPA startups snagged significant funding as well. Then again, the market is forecasted to see tremendous growth. Grand View Research predicts that spending will hit $3.97 billion in the United States by 2025.[4]

Interestingly enough, Forrester had this to say about the RPA trend:

> Today's most successful companies generally operate with fewer employees than those of the past. Consider that Kodak at its peak in 1973 employed 120,000, but when Facebook bought Instagram in 2012, the photo-sharing site employed only 13 workers. In 2019, we predict that one in 10 startups—operating in a more agile, lean, and scalable fashion— will look at the world through the lens of tasks, not jobs, and will build business models around automation-first principles.[5]

---

[2] http://business-review.eu/news/the-story-of-uipath-how-it-became-romanias-first-unicorn-164248

[3] www.uipath.com/newsroom/uipath-raises-225-million-series-c-led-by-capitalg-and-sequoia

[4] www.grandviewresearch.com/press-release/global-robotic-process-automation-rpa-market

[5] https://go.forrester.com/blogs/predictions-2019-automation-will-become-central-to-business-strategy-and-operations/

RPA is yet another area that has been supercharged with AI. If anything, it could be the gateway for many companies because the implementation usually does not take long or require heavy costs.

In this chapter, we'll take a look at RPA and see how it could be a critical driver for many companies.

# What Is RPA?

The term Robotic Process Automation can be a bit confusing. The word "robotic" does not mean physical robots (we'll cover these in Chapter 7); rather, it is about software-based robots or bots.

RPA allows you to use low-code visual drag-and-drop systems to automate the workflow of a process. Some examples include the following:

- Inputting, changing, and tracking Human Resources (HR) documents, contracts, and employee information

- Detecting issues with customer service and taking actions to resolve the problems

- Processing an insurance claim

- Sending invoices

- Issuing refunds to customers

- Reconciling financial records

- Transferring data from one system to another

- Providing standard replies to customers

This is all done by having a bot replicate the workflows for an application, say for an ERP (Enterprise Resource Planning) or CRM (Customer Relationship Management) system. This may even be done with the RPA program recording the steps from employees or with the use of OCR (optical character recognition) technology to translate handwritten notes. Think of RPA as a digital employee.

There are two flavors of this type of technology:

- *Unattended RPA*: This is a process that is completely autonomous as the bot will run in the background. Now this does not mean there is no human intervention. There will still be intervention for exception management. This is when the bot encounters something it does not understand.

- *RDA (Robotic Desktop Automation)*: This is where RPA helps an employee with a job or task. A common use case is with a contact center. That is, when a call comes in, the rep can use RDA to help find answers, send messages, pull customer profile information, and get insight on what to do next. The technology helps improve or augment the efficiency of the worker.

# Pros and Cons of RPA

Of course, quite a bit of time for a typical employee—in the back office—is spent on routine tasks. But with RPA, companies can often get a strong ROI (Return on Investment)—so long as the implementation is done right.

Here are some other advantages:

- *Customer Satisfaction*: RPA means minimal errors as well as high speed. A bot also works 24/7. This means that customer satisfaction scores—like the NPS (Net Promoter Score)—should improve. Note that increasingly more customers, such as from the Millennial generation, prefer to deal with apps/web sites, not people! RPA also means that reps will have more time to spend on value-added tasks, instead of dealing with the tedious matters that waste time.

- *Scalability*: Once a bot is created, it can be quickly expanded to meet spikes in activity. This can be critical for seasonal businesses like retailers.

- *Compliance*: For people, it's tough to keep track of rules, regulations, and laws. Even worse, they often change. But with RPA, compliance is built into the process—and is always followed. This can be a major benefit in terms of avoiding legal problems and fines.

- *Insights and Analytics*: Next-generation RPA platforms come equipped with sophisticated dashboards, which focus on KPIs for your business. You can also set up alerts if there are any problems.

- *Legacy Systems*: Older companies are often bogged down with old IT systems, which makes it extremely tough to pull off a digital transformation. But RPA software is able to work fairly well with legacy IT environments.

- *Data*: Because of the automation, the data is much cleaner as there are minimal input errors. This means that organizations will—over time—have more accurate understandings of their businesses. The data quality will also increase the likelihood of success of AI implementations.

While all this is great, RPA still has its drawbacks. For example, if you have current processes that are inefficient and you rush to implement the RPA system, you will be essentially replicating a bad approach! This is why it is critical to evaluate your workflows before implementing a system.

But there are certainly other potential landmines to note, such as the following:

- *Brittleness*: RPA can easily break if there are changes in the underlying applications. This could also be the case if there are changes in procedures and regulations. It's true that newer systems are getting better at adapting and may also leverage APIs. But RPA is not about hands-off activity.

- *Virtualized Apps*: This type of software, such as from Citrix, can be difficult with RPA systems because they cannot effectively capture the processes. The reason is that the data is stored on an outside server and the output is a snapshot on a monitor. But some companies are using AI to solve the problem, such as UiPath. The company has created a system, called "Pragmatic AI," which uses computer vision to interpret the screen snapshots to record the processes.

- *Specialization*: Many RPA tools are for general-purpose activities. But there may be areas that require specialization, such as with finance. In this case, you may look at a niche software app that can handle it.

- *Testing*: This is absolutely critical. You want to first sample some transactions to make sure the system is working correctly. After this, you can do a more extensive rollout of the RPA system.

- *Ownership*: The temptation is to have IT own the RPA implementation and management. But this is probably not advisable. The reason? RPA systems are fairly low tech. After all, they can be developed by nonprogrammers. Because of this, the business managers are ideal for owning the process since they can generally handle the technical issues and also have a firmer grasp of the employee workflows.

- *Resistance*: Change is always difficult. With RPA, there may be fears that the technology will displace jobs. This means you need to have a clear set of messages, which focus on the benefits of the technology. For example, RPA will mean more time to focus on important matters, which should make a person's job more interesting and meaningful.

# What Can You Expect from RPA?

When it comes to RPA, the industry is still in the early phases. Yet there are clear signs that the technology is making a big difference for many companies.

Take a look at the research report from Computer Economics Technology, which included roughly 250 companies (they were across many industries and had revenues that were from $20 million to over $1 billion). Of those that implemented an RPA system, about half reported a positive return within 18 months of deployment. This is definitely standout for enterprise software, which can be challenging in getting adoption.[6]

And to get a sense of the strategic importance of this technology, look to see what the US Defense Department—which is engaged in over 500 AI projects—is doing. Here's what the agency's Joint Artificial Intelligence Center director, Air Force Lt. Gen. Jack Shanahan, had to say during a Congressional hearing:

> When you talk about smart automation, or in the vernacular of the industry, Robotic Process Automation, it's not headline grabbing in terms of big AI projects, but it may be where the most efficiencies can be found. That's the case if you read some of the dailies in industry, whether it's in medicine or finance, this is where early gains are being realized in AI. Some of the other projects we take on in the department are probably years in the making in return on investment.[7]

Despite all this, there are still many failed RPA implementations as well. Ernst & Young, for example, has received large amount of consulting business because of this. Based on this experience, the failure rate for initial RPA projects ranges from 30% to 50%.[8]

---

[6] www.computereconomics.com/article.cfm?id=2633
[7] https://federalnewsnetwork.com/artificial-intelligence/2019/03/dod-laying-groundwork-for-multi-generational-effort-on-ai/
[8] www.cmswire.com/information-management/why-rpa-implementation-projects-fail/

But this is inevitable with any type of enterprise software system. Yet so far, the problems appear mostly to be about planning, strategy, and expectations—not the technology.

Another problem is that the hype surrounding RPA may be raising expectations to excessive levels. This means that disappointment will be fairly common, even if implementations are successful!

Of course, technologies are not cure-alls. And they require much time, effort, and diligence to work.

# How to Implement RPA

Then what are some steps to take for a successful RPA implementation? There is no standard answer, but there are certainly some best practices emerging:

- Determine the right functions to automate.
- Assess the processes.
- Select the RPA vendor and deploy the software.
- Set in place a team to manage the RPA platform.

Let's take a closer look at each of these.

# Determine the Right Functions to Automate

*Yes, excessive automation at Tesla was a mistake. To be precise, my mistake. Humans are underrated.*

—Elon Musk, CEO of Tesla[9]

Even though RPA is powerful and can move the needle in a big way for a company, the capabilities are still fairly limited. The technology essentially makes the most sense for automating repetitive, structured, and routine processes. This involves things like scheduling, inputting/transferring data, following rules/workflows, cut and paste, filling out forms, and search. This means that RPA can actually have a role in just about every department in an organization.

Then where does this technology generally fail to deliver? Well, if a process requires independent judgment, then RPA probably does not make sense. The same goes for when the processes are subject to frequent change. In this situation, you can spend lots of time with ongoing adjustments to the configurations.

---

[9] https://twitter.com/elonmusk/status/984882630947753984?lang=en

Once you establish a part of the business where the technology looks like a good fit, there are a myriad of other considerations. In other words, you'll likely have more success with a project if you focus on the following:

- The areas of the business that have serious levels of underperformance

- The processes that take up a high percentage of employee time and involve high error rates

- The tasks that need more hiring when there are higher volumes

- The areas that you are thinking of outsourcing

- A process that has a large number of steps and in which there are various applications involved

## Assess the Processes

It's common for a company to have many unwritten processes. And this is fine. This approach allows for adaptability, which is what people are good at.

However, this is far from the case with a bot. To have a successful implementation, you need to have a deep assessment of the processes. This can actually take a while, and it may make sense to get outside consultants to help out. They have the advantage of being more neutral and better able to identify the weaknesses.

Some of the RPA vendors do have their own tools to help with analyzing processes—which you should definitely use. There are also third-party software providers that have their own offerings. One is Celonis, which integrates with RPA platforms such as UiPath, Automation Anywhere, Blue Prism, and others. The software performs essentially a digital MRI that analyzes data, providing insights on how your processes really work. It will also identify weakness and opportunities, such as to increase revenues, improve customer satisfaction, and free up resources.

Regardless of what approach you take, it is critical that you formulate a clear-cut plan that has the input from IT, higher management, and the departments impacted. Also make sure to get analytics people involved, as there could be opportunities for leveraging the data.

## Select the RPA Vendor and Deploy the Software

By going through the first two steps, you'll be in a very good position to evaluate the different RPA systems. For example, if your main goal is to cut staff, then you would look for software that is focused on unattended bots.

Or, if you want to leverage data—such as for AI applications—then this will lead to other types of RPA platforms.

Then, once you have selected one, you will begin the deployment. The good news is that this can be relatively fast, say less than a month.

But as you go on to do further RPA projects, you may run into something called automation fatigue. This is where the returns generally start to deteriorate.

Think of it this way: When you start out, the focus will usually be on those areas of the business that need automation the most, which means the ROI will be significant. But over time, there will be a focus on tasks that are not as amenable to automation, and it will likely take much more work to realize even slight improvements.

Because of this, it is a good idea to temper expectations when engaging in a widespread RPA transformation.

## Set in Place a Team to Manage the RPA Platform

Just because RPA provides a high degree of automation does not mean it requires little management. Rather, the best approach is to put together a team, which is often referred to as a center of excellence (CoE).

In order to make the best use of the CoE, you need to be clear on each person's responsibilities. For example, you should be able to answer the following questions:

- What happens if there is a problem with a bot? At what points should there be human intervention?
- Who is in charge of monitoring the RPA?
- Who is in charge of training?
- Who will have the role for the first line of support?
- Who is responsible for developing the bots?

For larger organizations, you might also want to expand the roles. You could have an RPA champion, who would be the evangelist of the platform—for the whole company. Or there could be an RPA change manager, who provides the communication to help with adoption.

Finally, as the RPA implementation gets larger, a key goal should be to look at how all the parts fit together. Like many other software systems, there is the risk of sprawl across the organization—which can mean not getting higher performance. This is where having a proactive CoE can make a major positive impact.

# RPA and AI

While still in the initial phases, AI is already making strides with RPA tools. This is leading to the emergence of Cognitive Robotic Process Automation (CRPA) software bots.

And this makes sense. After all, RPA is about optimizing processes and involves large amounts of data. So vendors are starting to implement systems like machine learning, deep learning, speech recognition, and Natural Language Processing. Some of the leaders in the CRPA space include UiPath, Automation Anywhere, Blue Prism, NICE Systems, and Kryon Systems.

For example, with Automation Anywhere, a bot can handle tasks such as extracting invoices from emails, which involves sophisticated text processing. The company also has prebuilt integrations with third-party AI services like IBM Watson, AWS Machine Learning, and Google Cloud AI.[10]

"There has been a proliferation of AI-enabled services in recent years, but businesses often struggle to operationalize them," said Mukund Srigopal, who is the director of Product Marketing at Automation Anywhere. "RPA is a great way to infuse AI capabilities into business processes."[11]

Here are some other ways CRPA can allow for AI functions:

- You can connect chatbots with your system, which will allow for automated customer service (we'll cover this topic in Chapter 6).

- AI can find the right moment to send an email or alert.

- IVR (Interactive Voice Response) has gotten a bad reputation over the years. Simply put, customers do not like the hassle of going through multiple steps to solve a problem. But with CRPA, you can use something called Dynamic IVR. This personalizes the voice messages to each customer, providing for a much better experience.

- NLP and text analysis can convert unstructured data into structured data. This can make the CRPA more effective.

---

[10] www.forbes.com/sites/tomtaulli/2019/02/02/what-you-need-to-know-about-rpa-robotic-process-automation/
[11] This is from the author's interview with Mukund Srigopal, the director of Product Marketing at Automation Anywhere.

# RPA in the Real World

To get a better sense of how RPA works and to understand the benefits, here's a look at a case study of Microsoft.[12] Every year, the company pays billions of dollars in royalties to game developers, partners, and content creators. Yet the process was mostly manual, involving the sending of thousands of statements—and yes, this was a big time waster for the company.

So the company selected Kyron for an RPA implementation. By doing an initial process review, Microsoft realized that anywhere from 70% to 80% of the statements were straightforward and could be easily automated. The rest included exceptions that required human intervention, such as approvals.

With the RPA system, a visual detection algorithm could divvy up the statements and find the exceptions. The setup was also fairly quick, taking about 6 weeks.

As should be no surprise, the results had a material impact on the process. For instance, a bot was able to take only 2.5 hours to complete 150 royalty statements. By comparison, it would take 50 hours for employees to do this. The bottom line: Microsoft achieved a 2,000% savings. There was also an elimination of any rework from human error (which before was about 5% in a given month).

# Conclusion

As seen with the Microsoft case study, RPA can lead to major savings. But there still needs to be diligent planning, so as to understand your processes. For the most part, the focus should be on tasks that are manual and repetitive—not those that rely heavily on judgment. Next, it is important to setup a CoE to oversee the ongoing management of the automation, which will help with handling exceptions, collecting data, and tracking KPIs.

RPA is also a great way to implement basic AI within an organization. Actually, because there could be significant ROI, this may spur even more investment in pursuing this technology.

---

[12]www.kryonsystems.com/microsoft-case-study/

# Key Takeaways

- Robotic Process Automation (RPA) allows you to use low-code visual drag-and-drop systems to automate the workflow of a process.

- Unattended RPA is when a process is completely automated.

- RDA (Robotic Desktop Automation) is where RPA helps an employee with a job or task.

- Some of the benefits of RPA include higher customer satisfaction, lower error rates, improved compliance, and easier integration with legacy systems.

- Some of the drawbacks of RPA include the difficulty with adapting to changes in the underlying applications, issues with virtualized apps, and resistance from employees.

- RPA tends to work best where you can automate repetitive, structured, and routine processes, such as scheduling, inputting/transferring data, and following rules/workflows.

- When implementing an RPA solution, some of the steps to consider include determining the functions to automate, assessing the processes, selecting the RPA vendor and deploying the software, and setting in place a team to manage the platform.

- A center of excellence (CoE) is a team that manages an RPA implementation.

- Cognitive Robotic Process Automation (CRPA) is an emerging category of RPA that focuses on AI technologies.

# Natural Language Processing (NLP)

## How Computers Talk

In 2014, Microsoft launched a chatbot—an AI system that communicates with people—called Xiaoice. It was integrated into Tencent's WeChat, the largest social messaging service in China. Xiaoice performed quite well, getting to 40 million users within a few years.

In light of the success, Microsoft wanted to see if it could do something similar in the US market. The company's Bing and the Technology and Research Group leveraged AI technologies to build a new chatbot: Tay. The developers even enlisted the help of improvisational comedians to make the conversion engaging and fun.

On March 23, 2016, Microsoft launched Tay on Twitter—and it was an unmitigated disaster. The chatbot quickly spewed racist and sexist messages! Here's just one of the thousands of examples:

> *@TheBigBrebowski ricky gervais learned totalitarianism from adolf hitler, the inventor of atheism*[1]

Tay was a vivid illustration of Godwin's Law. It reads as follows: the more an online discussion continues, the higher are the odds that someone will bring up Adolf Hitler or the Nazis.

So yes, Microsoft took down Tay within 24 hours and blogged an apology. In it, the corporate vice president of Microsoft Healthcare, Peter Lee, wrote:

> Looking ahead, we face some difficult—and yet exciting— research challenges in AI design. AI systems feed off of both positive and negative interactions with people. In that sense, the challenges are just as much social as they are technical. We will do everything possible to limit technical exploits but also know we cannot fully predict all possible human interactive misuses without learning from mistakes. To do AI right, one needs to iterate with many people and often in public forums. We must enter each one with great caution and ultimately learn and improve, step by step, and to do this without offending people in the process. We will remain steadfast in our efforts to learn from this and other experiences as we work toward contributing to an Internet that represents the best, not the worst, of humanity.[2]

A key to Tay was to repeat some of the content of the people asking questions. For the most part, this is a valid approach. As we saw in Chapter 1, this was at the heart of the first chatbot, ELIZA.

But there also must be effective filters in place. This is especially the case when a chatbot is used in a free-form platform like Twitter (or, for that matter, in any real-world scenario).

However, failures like Tay are important. They allow us to learn and to evolve the technology.

In this chapter, we'll take a look at chatbots as well as Natural Language Processing (NLP), which is a key part of how computers understand and manipulate language. This is a subset of AI.

Let's get started.

---

[1] www.theverge.com/2016/3/24/11297050/tay-microsoft-chatbot-racist
[2] https://blogs.microsoft.com/blog/2016/03/25/learning-tays-introduction/

# The Challenges of NLP

As we saw in Chapter 1, language is the key to the Turing Test, which is meant to validate AI. Language is also something that sets us apart from animals.

But this area of study is exceedingly complex. Here are just some of the challenges with NLP:

- Language can often be ambiguous. We learn to speak in a quick fashion and accentuate our meaning with nonverbal cues, our tone, or reactions to the environment. For example, if a golf ball is heading toward someone, you'll yell "Fore!" But an NLP system would likely not understand this because it cannot process the context of the situation.

- Language changes frequently as the world changes. According to the Oxford English Dictionary, there were more than 1,100 words, senses, and subentries in 2018 (in all, there are over 829,000)[3]. Some of the new entries included mansplain and hangry.

- When we talk, we make grammar mistakes. But this is usually not a problem as people have a great ability for inference. But this is a major challenge for NLP as words and phrases may have multiple meanings (this is called polysemy). For example, noted AI researcher Geoffrey Hinton likes to compare "recognize speech" and "wreck a nice beach."[4]

- Language has accents and dialects.

- The meaning of words can change based on, say, the use of sarcasm or other emotional responses.

- Words can be vague. After all, what does it really mean to be "late"?

- Many words have essentially the same meaning but involve degrees of nuances.

- Conversations can be non-linear and have interruptions.

Despite all this, there have been great strides with NLP, as seen with apps like Siri, Alexa, and Cortana. Much of the progress has also happened within the last decade, driven by the power of deep learning.

---

[3] https://wordcounter.io/blog/newest-words-added-to-the-dictionary-in-2018/
[4] www.deepinstinct.com/2019/04/16/applications-of-deep-learning/

Now there can be confusion about human languages and computer languages. Haven't computers been able to understand languages like BASIC, C, and C++ for years? This is definitely true. It's also true that computer languages have English words like *if*, *then*, *let*, and *print*.

But this type of language is very different from human language. Consider that a computer language has a limited set of commands and strict logic. If you use something incorrectly, this will result in a bug in the code—leading to a crash. Yes, computer languages are very literal!

# Understanding How AI Translates Language

Now as we saw in Chapter 1, NLP was an early focus for AI researchers. But because of the limited computer power, the capabilities were quite weak. The goal was to create rules to interpret words and sentences—which turned out to be complex and not very scalable. In a way, NLP in the early years was mostly like a computer language!

But over time, there evolved a general structure for it. This was critical since NLP deals with unstructured data, which can be unpredictable and difficult to interpret.

Here's a general high-level look at the two key steps:

- *Cleaning and Preprocessing the Text*: This involves using techniques like tokenization, stemming, and lemmatization to parse the text.

- *Language Understanding and Generation*: This is definitely the most intensive part of the process, which often uses deep learning algorithms.

In the next few sections, we'll look at the different steps in more detail.

## Step #1—Cleaning and Preprocessing

Three things need to be done during the cleaning and preprocessing step: tokenization, stemming, and lemmatization.

### Tokenization

Before there can be NLP, the text must be parsed and segmented into various parts—a process known as tokenization. For example, let's say we have the following sentence: "John ate four cupcakes." You would then separate and categorize each element. Figure 6-1 illustrates this tokenization.
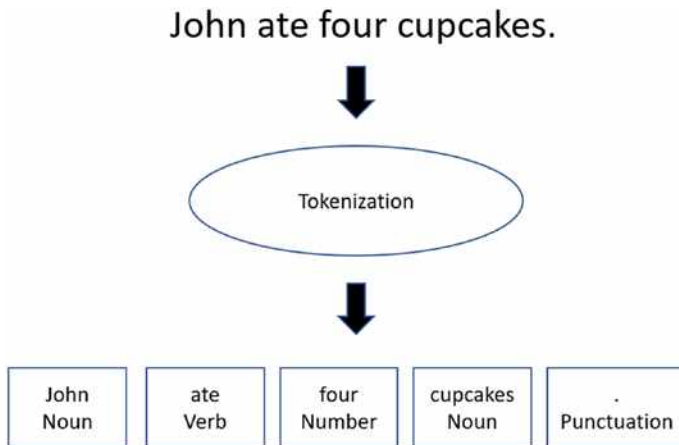
**Figure 6-1.** Example of sentence tokenization

All in all, kind of easy? Kind of.

After tokenization, there will be normalization of the text. This will entail converting some of the text so as to make it easier for analysis, such as by changing the case to upper or lower, removing punctuation, and eliminating contractions.

But this can easily lead to some problems. Suppose we have a sentence that has "A.I." Should we get rid of the periods? And if so, will the computer know what "A I" means?

Probably not.

Interestingly enough, even the case of words can have a major impact on the meaning. Just look at the difference between "fed" and the "Fed." The Fed is often another name for the Federal Reserve. Or, in another case, let's suppose we have "us" and "US." Are we talking about the United States here?

Here are some of the other issues:

- *White Space Problem*: This is where two or more words should be one token because the words form a compound phrase. Some examples include "New York" and "Silicon Valley."

- *Scientific Words and Phrases*: It's common for such words to have hyphens, parentheses, and Greek letters. If you strip out these characters, the system may not be able to understand the meanings of the words and phrases.

- *Messy Text*: Let's face it, many documents have grammar and spelling errors.

- *Sentence Splitting*: Words like "Mr." or "Mrs." can prematurely end a sentence because of the period.

- *Non-important Words*: There are ones that really add little or no meaning to a sentence, like "the," "a," and "an." To remove these, you can use a simple Stop Words filter.

As you can see, it can be easy to mis-parse sentences (and in some languages, like Chinese and Japanese, things can get even more difficult with the syntax). But this can have far-ranging consequences. Since tokenization is generally the first step, a couple errors can cascade through the whole NLP process.

## Stemming

Stemming describes the process of reducing a word to its root (or lemma), such as by removing affixes and suffixes. This has actually been effective for search engines, which involve the use of clustering to come up with more relevant results. With stemming, it's possible to find more matches as the word has a broader meaning and even to handle such things as spelling errors. And when using an AI application, it can help improve the overall understanding.
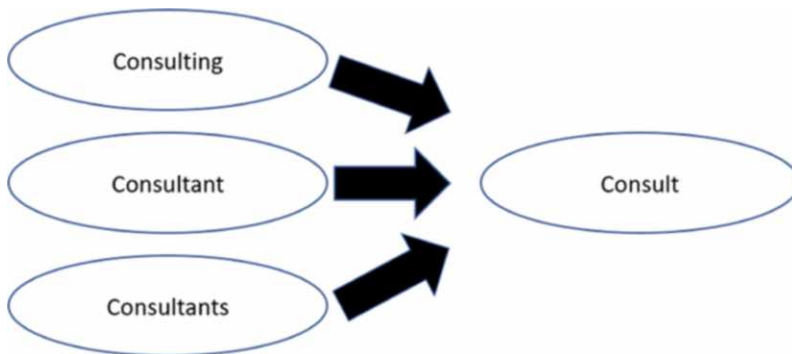
Figure 6-2 shows an example of stemming.



**Figure 6-2.** Example of stemming

There are a variety of algorithms to stem words, many of which are fairly simple. But they have mixed results. According to IBM:

> The Porter algorithm, for example, will state that 'universal' has the same stem as 'university' and 'universities,' an observation that may have historical basis but is no longer semantically relevant. The Porter stemmer also does not recognize that 'theater' and 'theatre' should belong to the same stem class. For reasons such as these, Watson Explorer Engine does not use the Porter stemmer as its English stemmer.[5]

In fact, IBM has created its own proprietary stemmer, and it allows for significant customization.

## Lemmatization

Lemmatization is similar to stemming. But instead of removing affixes or prefixes, there is a focus on finding similar root words. An example is "better," which we could lemmatize to "good." This works so long as the meaning remains mostly the same. In our example, both are roughly similar, but "good" has a clearer meaning. Lemmatization also may work with providing better searches or language understanding, especially with translations.

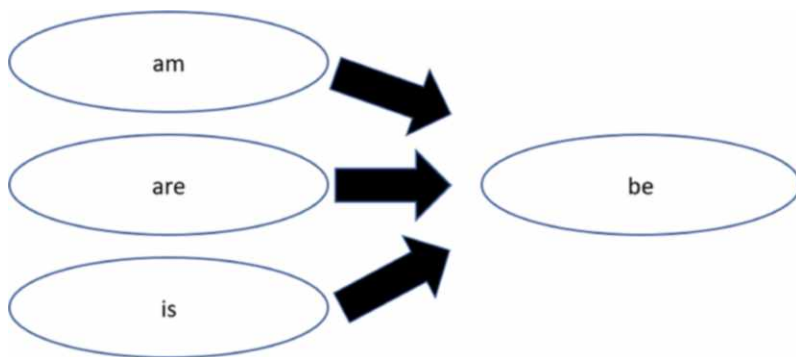Figure 6-3 shows an example of lemmatization.



**Figure 6-3.** Example of lemmatization

---

[5] `www.ibm.com/support/knowledgecenter/SS8NLW_11.0.1/com.ibm.swg.im.infosphere.dataexpl.engine.doc/c_correcting_stemming_errors.html`

To effectively use lemmatization, the NLP system must understand the meanings of the words and the context. In other words, this process usually has better performance than stemming. On the other hand, it also means that the algorithms are more complicated and there are higher levels of computing power required.

# Step #2—Understanding and Generating Language

Once the text has been put into a format that computers can process, then the NLP system must understand the overall meaning. For the most part, this is the hardest part.

But over the years, researchers have developed a myriad of techniques to help out, such as the following:

- *Tagging Parts of Speech (POS)*: This goes through the text and designates each word into its proper grammatical form, say nouns, verbs, adverbs, etc. Think of it like an automated version of your grade school English class! What's more, some POS systems have variations. Note that a noun has singular nouns (NN), singular proper nouns (NNP), and plural nouns (NNS).

- *Chunking*: The words will then be analyzed in terms of phrases. For example, a noun phrase (NP) is a noun that acts as the subject or object to a verb.

- *Named Entity Recognition*: This is identifying words that represent locations, persons, and organizations.

- *Topic Modelling*: This looks for hidden patterns and clusters in the text. One of the algorithms, called Latent Dirichlet Allocation (LDA), is based on unsupervised learning approaches. That is, there will be random topics assigned, and then the computer will iterate to find matches.

For many of these processes, we can use deep learning models. They can be extended to more areas of analysis—to allow for seamless language understanding and generation. This is a process known as distributional semantics.

With a convolutional neural network (CNN), which we learned about in Chapter 4, you can find clusters of words that are translated into a feature map. This has allowed for applications like language translation, speech recognition, sentiment analysis, and Q&A. In fact, the model can even do things like detect sarcasm!

Yet there are some problems with CNNs. For example, the model has difficulties with text that has dependencies across large distances. But there are some ways to handle this, such as with time-delayed neural networks (TDNN) and dynamic convolutional neural networks (DCNN). These methods have shown high performance in handling sequenced data. Although, the model that has shown more success with this is the recurrent neural network (RNN), as it memorizes data.

So far, we have been focused mostly on text analysis. But for there to be sophisticated NLP, we also must build voice recognition systems. We'll take a look at this in the next section.

# Voice Recognition

In 1952, Bell Labs created the first voice recognition system, called Audrey (for Automatic Digit Recognition). It was able to recognize phonemes, which are the most basic units of sounds in a language. English, for example, has 44.

Audrey could recognize the sound of a digit, from zero to nine. It was accurate for the voice of the machine's creator, HK Davis, about 90% of the time.[6] And for anyone else, it was 70% to 80% or so.

Audrey was a major feat, especially in light of the limited computing power and memory available at the time. But the program also highlighted the major challenges with voice recognition. When we speak, our sentences can be complex and somewhat jumbled. We also generally talk fast—an average of 150 words per minute.

As a result, voice recognition systems improved at a glacially slow pace. In 1962, IBM's Shoebox system could recognize only 16 words, 10 digits, and 6 mathematical commands.

It was not until the 1980s that there was significant progress in the technology. The key breakthrough was the use of the hidden Markov model (HMM), which was based on sophisticated statistics. For example, if you say the word "dog," there will be an analysis of the individual sounds d, o, and g. The HMM algorithm will assign a score to each of these. Over time, the system will get better at understanding the sounds and translate them into words.

While HMM was critical, it still was unable to effectively handle continuous speech. For example, voice systems were based on template matching. This involved translating sound waves into numbers, which was done by sampling. The result was that the software would measure the frequency of the intervals and store the results. But there had to be a close match. Because of this, the voice input had to be quite clear and slow. There also had to be little background noise.

---

[6] www.bbc.com/future/story/20170214-the-machines-that-learned-to-listen

But by the 1990s, software developers would make strides and come out with commercial systems, such as Dragon Dictate, which could understand thousands of words in continuous speech. However, adoption was still not mainstream. Many people still found it easier to type into their computers and use the mouse. Yet there were some professions, like medicine (a popular use case with transcribing diagnosis of patients), where speech recognition found high levels of usage.

With the emergence of machine learning and deep learning, voice systems have rapidly become much more sophisticated and accurate. Some of the key algorithms involve the use of the long short-term memory (LSTM), recurrent neural networks, and deep feed-forward neural networks. Google would go on to implement these approaches in Google Voice, which was available to hundreds of millions of smartphone users. And of course, we've seen great progress with other offerings like Siri, Alexa, and Cortana.

# NLP in the Real World

For the most part, we have gone through the main parts of the NLP workflow. Next, let's take a look at the powerful applications of this technology.

## Use Case: Improving Sales

Roy Raanani, who has a career in working with tech startups, thought that the countless conversions that occur every day in business are mostly ignored. Perhaps AI could transform this into an opportunity?

In 2015, he founded Chorus to use NLP to divine insights from conversations from sales people. Raanani called this the Conversation Cloud, which records, organizes, and transcribes calls—which are entered in a CRM (Customer Relationship Management) system. Over time, the algorithms will start to learn about best practices and indicate how things can be improved.

But pulling this off has not been easy. According to a Chorus blog:

> There are billions of ways to ask questions, raise objections, set action items, challenge hypotheses, etc. all of which need to be identified if sales patterns are to be codified. Second, signals and patterns evolve: new competitors, product names and features, and industry-related terminology change over time, and machine-learned models quickly become obsolete.[7]

---

[7] https://blog.chorus.ai/a-taste-of-chorus-s-secret-sauce-how-our-system-teaches-itself

For example, one of the difficulties—which can be easily overlooked—is how to identify the parties who are talking (there are often more than three on a call). Known as "speaker separation," it is considered even more difficult than speech recognition. Chorus has created a deep learning model that essentially creates a "voice fingerprint"—which is based on clustering—for each speaker. So after several years of R&D, the company was able to develop a system that could analyze large amounts of conversations.

As a testament to this, look at one of Chorus' customers, Housecall Pro, which is a startup that sells mobile technologies for field service management. Before adopting the software, the company would often create personalized sales pitches for each lead. But unfortunately, it was unscalable and had mixed results.

But with Chorus, the company was able to create an approach that did not have much variation. The software made it possible to measure every word and the impact on the sales conversions. Chorus also measured whether a sales rep was on-script or not.

The outcome? The company was able to increase the win rate of the sales organization by 10%.[8]

## Use Case: Fighting Depression

Across the world, about 300 million people suffer from depression, according to data from the World Health Organization.[9] About 15% of adults will experience some type of depression during their life.

This may go undiagnosed because of lack of healthcare services, which can mean that a person's situation could get much worse. Unfortunately, the depression can lead to other problems.

But NLP may be able to improve the situation. A recent study from Stanford used a machine learning model that processed 3D facial expressions and the spoken language. As a result, the system was able to diagnose depression with an average error rate of 3.67 when using the Patient Health Questionnaire (PHQ) scale. The accuracy was even higher for more aggravated forms of depression.

In the study, the researchers noted: "This technology could be deployed to cell phones worldwide and facilitate low-cost universal access to mental health care."[10]

---

[8] www.chorus.ai/case-studies/housecall/
[9] www.verywellmind.com/depression-statistics-everyone-should-know-4159056
[10] "*Measuring Depression Symptom Severity from Spoken Language and 3D Facial Expressions*," A Haque, M Guo, AS Miner, L Fei-Fei, presented at the NeurIPS 2018 Workshop on Machine Learning for Health (ML4H), https://arxiv.org/abs/1811.08592.

# Use Case: Content Creation

In 2015, several tech veterans like Elon Musk, Peter Thiel, Reid Hoffman, and Sam Altman launched OpenAI, with the support of a whopping $1 billion in funding. Structured as a nonprofit, the goal was to develop an organization with the goal "to advance digital intelligence in the way that is most likely to benefit humanity as a whole, unconstrained by a need to generate financial return."[11]

One of the areas of research has been on NLP. To this end, the company launched a model called GPT-2 in 2019, which was based on a dataset of roughly eight million web pages. The focus was to create a system that could predict the next word based on a group of text.

To illustrate this, OpenAI provided an experiment with the following text as the input: "In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English."

From this, the algorithms created a convincing story that was 377 words in length!

Granted, the researchers admitted that the storytelling was better for topics that related more to the underlying data, on topics like *Lord of the Rings* and even Brexit. As should be no surprise, GPT-2 demonstrated poor performance for technical domains.

But the model was able to score high on several well-known evaluations of reading comprehension. See Table 6-1.[12]

**Table 6-1.** Reading comprehension results

| DataSet | Prior Record for Accuracy | GPT-2's Accuracy |
|---|---|---|
| Winograd Schema Challenge | 63.7% | 70.70% |
| LAMBADA | 59.23% | 63.24% |
| Children's Book Test Common Nouns | 85.7% | 93.30% |
| Children's Book Test Named Entities | 82.3% | 89.05% |

[11] https://openai.com/blog/introducing-openai/
[12] https://openai.com/blog/better-language-models/

Even though a typical human would score 90%+ on these tests, the performance of GPT-2 is still impressive. It's important to note that the model used Google's neural network innovation, called a Transformer, and unsupervised learning.

In keeping with OpenAI's mission, the organization decided not to release the complete model. The fear was that it could lead to adverse consequences like fake news, spoofed Amazon.com reviews, spam, and phishing scams.

# Use Case: Body Language

Just focusing on language itself can be limiting. Body language is also something that should be included in a sophisticated AI model.

This is something that Rana el Kaliouby has been thinking about for some time. While growing up in Egypt, she earned her master's degree in science from the American University in Cairo and then went on to get her PhD in computer science at Newnham College of the University of Cambridge. But there was something that was very compelling to her: How can computers detect human emotions?

However, in her academic circles, there was little interest. The consensus view in the computer science community was that this topic was really not useful.

But Rana was undeterred and teamed up with noted professor Rosalind Picard to create innovative machine learning models (she wrote a pivotal book, called *Affective Computing*, which looked at emotions and machines).[13] Yet there also had to be the use of other domains like neuroscience and psychology. A big part of this was leveraging the pioneering work of Paul Ekman, who did extensive research on human emotions based on a person's facial muscles. He found that there were six universal human emotions (wrath, grossness, scaredness, joy, loneliness, and shock) that could be coded by 46 movements called action units—all becoming a part of the Facial Action Coding System, or FACS.

While at the MIT Media Lab, Rana developed an "emotional hearing aid," which was a wearable that allowed those people with autism to better interact in social environments.[14] The system would detect the emotions of people and provide appropriate ways to react.

---

[13] Rosalind W. Picard, *Affective Computing* (MIT Press).
[14] www.newyorker.com/magazine/2015/01/19/know-feel

It was groundbreaking as the *New York Times* named it as one of the most consequential innovations in 2006. But Rana's system also caught the attention of Madison Avenue. Simply put, the technology could be an effective tool to gauge an audience's mood about a television commercial.

Then a couple years later, Rana launched Affectiva. The company quickly grew and attracted substantial amounts of venture capital (in all, it has raised $54.2 million).

Rana, who was once ignored, had now become one of the leaders in a trend called "emotion-tracking AI."

The flagship product for Affectiva is Affdex, which is a cloud-based platform for testing audiences for video. About a quarter of the Fortune Global 500 use it.

But the company has developed another product, called Affectiva Automotive AI, which is an in-cabin sensing system for a vehicle. Some of the capabilities include the following:

- Monitoring for driver fatigue or distraction, which will trigger an alert (say a vibration of the seat belt).

- Providing for a handoff to a semi-autonomous system if the driver is not waking or is angry. There is even an ability to provide route alternatives to lessen the potential for road rage!

- Personalizing the content—say music—based on the passenger's emotions.

For all of these offerings, there are advanced deep learning systems that process enormous amounts of features of a database that has more than 7.5 million faces. These models also account for cultural influences and demographic differences—which is all done in real-time.

## Voice Commerce

NLP-driven technologies like virtual assistants, chatbots, and smart speakers are poised to have powerful business models—and may even disrupt markets like e-commerce and marketing. We have already seen an early version of this with Tencent's WeChat franchise. The company, which was founded during the heyday of the Internet boom in the late 1990s, started with a simple PC-based messenger product called OICQ. But it was the introduction of WeChat that was a game changer, which has since become China's largest social media platform with over 1 billion monthly active users.[15]

---

[15] www.wsj.com/articles/iphones-toughest-rival-in-china-is-wechat-a-messaging-app-1501412406

But this app is more than for exchanging messages and posting content. WeChat has quickly morphed into an all-purpose virtual assistant, where you can easily hail a ride-sharing service, make a payment at a local retailer, place a reservation for a flight, or play a game. For example, the app accounts for close to 35% of the entire usage time on smartphones in China on a monthly basis. WeChat is also a major reason that the country has become increasingly a cash-less society.

All this points to the power of an emerging category called voice commerce (or v-commerce), where you can make purchases via chat or voice. It's such a critical trend that Facebook's Mark Zuckerberg wrote a blog post,[16] in early 2019, where he said the company would become more like…WeChat.

According to research from Juniper, the market for voice commerce is forecasted to hit a whopping $80 billion by 2023.[17] But in terms of the winners in this market, it seems like a good bet that it will be those companies that have large install bases of smart devices like Amazon, Apple, and Google. But there will still be room for providers of next-generation NLP technologies.

OK then, how might these AI systems impact the marketing industry? Well, to see how, there was an article in the *Harvard Business Review*, called "Marketing in the Age of Alexa" by Niraj Dawar and Neil Bendle. In it, the authors note that "AI assistants will transform how companies connect with their customers. They'll become the primary channel through which people get information, goods and services, and marketing will turn into the batter for their attention."[18]

Thus, the growth in chatbots, digital assistants, and smart speakers could be much bigger than the initial web-based e-commerce revolution. These technologies have significant benefits for customers, such as convenience. It's easy to tell a device to buy something, and the machine will also learn about your habits. So the next time you say you want to have a soft drink, the computer will know what you are referring to.

But this may lead to a winners-take-all scenario. Ultimately, it seems like consumers will use only one smart device for their shopping. In addition, for brands that want to sell their goods, there will be a need to deeply understand what customers really want, so as to become the preferred vendor within the recommendation engine.

---

[16] www.facebook.com/notes/mark-zuckerberg/a-privacy-focused-vision-for-social-networking/10156700570096634/
[17] https://voicebot.ai/2019/02/19/juniper-forecasts-80-billion-in-voice-commerce-in-2023-or-10-per-assistant/
[18] https://hbr.org/2018/05/marketing-in-the-age-of-alexa

# Virtual Assistants

In 2003, as the United States was embroiled in wars in the Middle East, the Defense Department was looking to invest in next-generation technologies for the battlefield. One of the key initiatives was to build a sophisticated virtual assistant, which could recognize spoken instructions. The Defense Department budgeted $150 million for this and tasked the SRI (Stanford Research Institute) Lab—based in Silicon Valley—to develop the application.[19] Even though the lab was a nonprofit, it still was allowed to license its technologies (like the inkjet printer) to startups.

And this is what happened with the virtual assistant. Some of the members of SRI—Dag Kittlaus, Tom Gruber, and Adam Cheyer—called it Siri and started their own company to capitalize on the opportunity. They founded the operation in 2007, which was when Apple's iPhone was launched.

But there had to be much more R&D to get the product to the point where it could be useful for consumers. The founders had to develop a system to handle real-time data, build a search engine for geographic information, and build security for credit cards and personal data. But it was NLP that was the toughest challenge.

In an interview, Cheyer noted:

> The hardest technical challenge with Siri was dealing with the massive amounts of ambiguity present in human language. Consider the phrase 'book 4-star restaurant in Boston'—seems very straightforward to understand. Our prototype system could handle this easily. However, when we loaded in tens of millions of business names and hundreds of thousands of cities into the system as vocabulary (just about every word in the English language is a business name), the number of candidate interpretations went through the roof.[20]

But the team was able to solve the problems and turn Siri into a powerful system, which was launched on Apple's App Store in February 2010. "It's the most sophisticated voice recognition to appear on a smartphone yet," according to a review in Wired.com.[21]

Steve Jobs took notice and called the founders. Within a few days, they would meet, and the discussions quickly led to an acquisition, which happened in late April for more than $200 million.

However, Jobs thought there needed to be improvements to Siri. Because of this, there was a re-release in 2011. This actually happened a day before Jobs died.

---

[19] www.huffingtonpost.com/2013/01/22/siri-do-engine-apple-iphone_n_2499165.html
[20] https://medium.com/swlh/the-story-behind-siri-fbeb109938b0
[21] www.wired.com/2010/02/siri-voice-recognition-iphone/

Fast forward to today, Siri has the largest market share position in the virtual assistant market, with 48.6%. Google Assistant is at 28.7%, and Amazon.com's Alexa has 13.2%.[22]

According to the "Voice Assistant Consumer Adoption Report," about 146.6 million people in the United States have tried virtual assistants on their smartphones and over 50 million with smart speakers. But this only covers part of the story. Voice technology is also becoming embedded into wearables, headphones, and appliances.[23]

Here are some other interesting findings:

- Using voice to search for products outranked searches for various entertainment options.

- When it comes to productivity, the most common use cases for voice include making calls, sending emails, and setting alarms.

- The most common use of voice on smartphones occurs when a person is driving.

- Regarding complaints with voice assistants on smartphones, the one with the highest percentage was inconsistency in understanding requests. Again, this points to the continuing challenges of NLP.

The growth potential for virtual assistants remains bright, and the category is likely to be a key for the AI industry. Juniper Research forecasts that the number of virtual assistants in use on a global basis will more than triple to 2.5 billion by 2023.[24] The fastest category is actually expected to be smart TVs. Yes, I guess we'll be holding conversations with these devices!

# Chatbots

There is often confusion between the differences between virtual assistants and chatbots. Keep in mind that there is much overlap between the two. Both use NLP to interpret language and perform tasks.

But there are still critical distinctions. For the most part, chatbots are focused primarily for businesses, such as for customer support or sales functions.

---

[22] www.businessinsider.com/siri-google-assistant-voice-market-share-charts-2018-6
[23] https://voicebot.ai/wp-content/uploads/2018/11/voice-assistant-consumer-adoption-report-2018-voicebot.pdf
[24] https://techcrunch.com/2019/02/12/report-voice-assistants-in-use-to-triple-to-8-billion-by-2023/

Virtual assistants, on the other hand, are geared for essentially everyone to help with their daily activities.

As we saw in Chapter 1, the origins of chatbots go back to the 1960s with the development of ELIZA. But it was not until the past decade or so that this technology became useable at scale.

Here's a sampling of interesting chatbots:

- *Ushur*: This is integrated in the enterprise systems for insurance companies, allowing for the automation of claims/bill processing and sales enablement. The software has shown, on average, a reduction of 30% in service center call volumes and a 90% customer response rate.[25] The company built its own state-of-the-art linguistics engine called LISA (this stands for Language Intelligence Services Architecture). LISA includes NLP, NLU, sentiment analysis, sarcasm detection, topic detection, data extraction, and language translations. The technology currently supports 60 languages, making it a useful platform for global organizations.

- *Mya*: This is a chatbot that can engage in conversations in the recruiting process. Like Ushur, this is also based on a home-grown NLP technology. Some of the reasons for this include having better communications but also handling specific topics for hiring.[26] Mya greatly reduces time to interview and time to hire by eliminating major bottlenecks.

- *Jane.ai*: This is a platform that mines data across an organization's applications and databases—say Salesforce.com, Office, Slack, and Gmail—in order to make it much easier to get answers, which are personalized. Note that about 35% of an employee's time is spent just trying to find information! For example, a use case of Jane.ai is USA Mortgage. The company used the technology, which was integrated into Slack, to help brokers to look up information for mortgage processing. The result is that USA Mortgage has saved about 1,000 human labor hours per month.[27]

---

[25] The information came from the author's interview with the CEO and co-founder of Ushur, Simha Sadasiva.

[26] The information came from the author's interview with the CEO and co-founder of Mya, Eyal Grayevsky.

[27] The information came from the author's interview with the CEO and co-founder of Jane.ai, David Karandish.

Despite all this, chatbots have still had mixed results. For example, just one of the problems is that it is difficult to program systems for specialized domains.

Take a look at a study from UserTesting, which was based on the responses from 500 consumers of healthcare chatbots. Some of the main takeaways included: there remains lots of anxiety with chatbots, especially when handling personal information, and the technology has problems with understanding complex topics.[28]

So before deploying a chatbot, there are some factors to consider:

- *Set Expectations*: Do not overpromise with the capabilities with chatbots. This will only set up your organization for disappointment. For example, you should not pretend that the chatbot is a human. This is a surefire way to create bad experiences. As a result, you might want to start off a chatbot conversation with "Hi, I'm a chatbot here to help you with…"

- *Automation*: In some cases, a chatbot can handle the whole process with a customer. But you should still have people in the loop. "The goal for chatbots is not to replace humans entirely, but to be the first line of defense, so to speak," said Antonio Cangiano, who is an AI evangelist at IBM. "This can mean not only saving companies money but also freeing up human agents who'll be able to spend more time on complex inquiries that are escalated to them."[29]

- *Friction*: As much as possible, try to find ways for the chatbot to solve problems as quickly as possible. And this may not necessarily be using a conversation. Instead, providing a simple form to fill out could be a better alternative, say to schedule a demo.

- *Repetitive Processes*: These are often ideal for chatbots. Examples include authentication, order status, scheduling, and simple change requests.

- *Centralization*: Make sure you integrate the data with your chatbots. This will allow for more seamless experiences. No doubt, customers quickly get annoyed if they have to repeat information.

---

[28] www.forbes.com/sites/bernardmarr/2019/02/11/7-amazing-examples-of-online-chatbots-and-virtual-digital-assistants-in-practice/#32bb1084533e
[29] This is from the author's interview with Antonio Cangiano, who is an AI evangelist at IBM.

- *Personalize the Experience*: This is not easy but can yield major benefits. Jonathan Taylor, who is the CTO of Zoovu, has this example: "Purchasing a camera lens will be different for every shopper. There are many variations of lenses that perhaps a slightly informed shopper understands—but the average consumer may not be as informed. Providing an assistive chatbot to guide a customer to the right lens can help provide the same level of customer service as an in-store employee. The assistive chatbot can ask the right questions, understanding the goal of the customer to provide a personalized product recommendation including 'what kind of camera do you already have,' 'why are you buying a new camera,' and 'what are you primarily trying to capture in your photographs?'"[30]

- *Data Analytics*: It's critical to monitor the feedback with a chatbot. What's the satisfaction? What's the accuracy rate?

- *Conversational Design and User Experience (UX)*: It's different than creating a web site or even a mobile app. With a chatbot, you need to think about the user's personality, gender, and even cultural context. Moreover, you must consider the "voice" of your company. "Rather than creating mockups of a visual interface, think about writing scripts and playing them out before to build it," said Gillian McCann, who is head of Cloud Engineering and Artificial Intelligence at Workgrid Software.[31]

Even with the issues with chatbots, the technology is continuing to improve. More importantly, these systems are likely to become an increasingly important part of the AI industry. According to IDC, about $4.5 billion will be spent on chatbots in 2019—which compares to a total of $35.8 billion estimated for AI systems.[32]

Something else: A study from Juniper Research indicates that the cost savings from chatbots are likely to be substantial. The firm predicts they will reach $7.3 billion by 2023, up from a mere $209 million in 2019.[33]

---

[30] This is from the author's interview with Jonathan Taylor, who is the CTO of Zoovu.
[31] This is from the author's interview with Gillian McCann, who is the head of Cloud Engineering and Artificial Intelligence at Workgrid Software.
[32] www.twice.com/retailing/artificial-intelligence-retail-chatbots-idc-spending
[33] www.juniperresearch.com/press/press-releases/bank-cost-savings-via-chatbots-to-reach

# Future of NLP

In 1947, Boris Katz was born in Moldova, which was part of the Soviet Union. He would go on to graduate from Moscow State University, where he learned about computers, and then left the country to the United States (with the assistance of Senator Edward Kennedy).

He wasted little time with the opportunity. Besides writing more than 80 technical publications and receiving two US patents, he created the START system that allowed for sophisticated NLP capabilities. It was actually the basis for the first Q&A site on the Web in 1993. Yes, this was the forerunner to breakout companies like Yahoo! and Google.

Boris's innovations were also critical for IBM's Watson, which is now at the core of the company's AI efforts. This computer, in 2011, would shock the world when it beat two of the all-time champions of the popular game show *Jeopardy!*

Despite all the progress with NLP, Boris is not satisfied. He believes we are still in the early stages and lots more must be done to get true value. In an interview with the *MIT Technology Review*, he said: "But on the other hand, these programs [like Siri and Alexa] are so incredibly stupid. So there's a feeling of being proud and being almost embarrassed. You launch something that people feel is intelligent, but it's not even close."[34]

This is not to imply he's a pessimist. However, he still thinks there needs to be a rethinking of NLP if it is to get to the point of "real intelligence." To this end, he believes researchers must look beyond pure computer science to broad areas like neuroscience, cognitive science, and psychology. He also thinks NLP systems must do a much better job of understanding the actions in the real world.

# Conclusion

For many people, the first interaction with NLP is with virtual assistants. Even while the technology is far from perfect, it still is quite useful—especially for answering questions or getting information, say about a nearby restaurant.

But NLP is also having a major impact in the business world. In the years ahead, the technology will become increasingly important for e-commerce and customer service—providing significant cost savings and allowing employees to focus on more value-added activities.

---

[34] www.technologyreview.com/s/612826/virtual-assistants-thinks-theyre-doomed-without-a-new-ai-approach/

True, there is still a long way to go because of the complexities of language. But the progress continues to be rapid, especially with the help of next-generation AI approaches like deep learning.

# Key Takeaways

- Natural Language Processing (NLP) is the use of AI to allow computers to understand people.

- A chatbot is an AI system that communicates with people, say by voice or online chat.

- While there have been great strides in NLP, there is much work to be done. Just some of the challenges include ambiguity of language, nonverbal cues, different dialects and accents, and changes to the language.

- The two main steps with NLP include cleaning/preprocessing the text and using AI to understand and generate language.

- Tokenization is where text is parsed and segmented into various parts.

- With normalization, text is converted into a form that makes it easier for analysis, such as by removing punctuation or contractions.

- Stemming describes the process of reducing a word to its root (or lemma), such as by removing affixes and suffixes.

- Similar to stemming, lemmatization involves finding similar root words.

- For NLP to understand language, there are a variety of approaches like tagging parts of speech (putting the text in the grammatical form), chunking (processing text in phrases), and topic modelling (finding hidden patterns and clusters).

- A phoneme is the most basic unit of sound in a language.

# Physical Robots

## The Ultimate Manifestation of AI

In the city of Pasadena, I went to CaliBurger for lunch and noticed a crowd of people next to the area where the food was being cooked—which was behind glass. The people were taking photos with their smartphones!

Why? The reason was Flippy, an AI-powered robot that can cook burgers.

I was there at the restaurant with David Zito, the CEO and co-founder of the company Miso Robotics that built the system. "Flippy helps improve the quality of the food because of the consistency and reduces production costs," he said. "We also built the robot to be in strict compliance with regulatory standards."[1]

After lunch, I walked over to the lab for Miso Robotics, which included a testing center with sample robots. It was here that I saw the convergence of software AI systems and physical robots. The engineers were building Flippy's brain, which was uploaded to the cloud. Just some of the capabilities included washing down utensils and the grill, learning to adapt to problems with cooking, switching between a spatula for raw meat and one for cooked meat, and placing baskets in the fryer. All this was being done in real-time.

But the food service industry is just one of the many areas that will be greatly impacted by robotics and AI.

---

[1] This is based on the author's interview, in January 2019, with David Zito, who is the CEO and co-founder of Miso Robotics.

According to International Data Corporation (IDC), the spending on robotics and drones is forecasted to go from $115.7 billion in 2019 to $210.3 billion by 2022.[2] This represents a compound annual growth rate of 20.2%. About two thirds of the spending will be for hardware systems.

In this chapter, we'll take a look at physical robots and how AI will transform the industry.

# What Is a Robot?

The origins of the word "robot" go back to 1921 in a play by Karel Capek called *Rossum's Universal Robots*. It's about a factory that created robots from organic matter, and yes, they were hostile! They would eventually join together to rebel against their human masters (consider that "robot" comes from the Czech word *robata* for forced labor).

But as of today, what is a good definition for this type of system? Keep in mind that there are many variations, as robots can have a myriad of forms and functions.

But we can boil them down into a few key parts:

- *Physical*: A robot can range in size, from tiny machines that can explore our body to massive industrial systems to flying machines to underwater vessels. There also needs to be some type of energy source, like a battery, electricity, or solar.

- *Act*: Simply enough, a robot must be able to take certain actions. This could include moving an item or even talking.

- *Sense*: In order to act, a robot must understand its environment. This is possible with sensors and feedback systems.

- *Intelligence*: This does not mean full-on AI capabilities. Yet a robot needs to be able to be programmed to take actions.

Nowadays it's not too difficult to create a robot from scratch. For example, RobotShop.com has hundreds of kits that range from under $10 to as much as $35,750.00 (this is the Dr. Robot Jaguar V6 Tracked Mobile Platform).

A heart-warming story of the ingenuity of building robots concerns a 2-year old, Cillian Jackson. He was born with a rare genetic condition that rendered him immobile. His parents tried to get reimbursement for a special electric wheelchair but were denied.

---

[2] www.idc.com/getdoc.jsp?containerId=prUS44505618

Well, the students at Farmington High School took action and built a system for Cillian.[3] Essentially, it was a robot wheelchair, and it took only a month to finish. Because of this, Cillian can now chase around his two corgis around the house!

While above we looked at the features of robots, there are key interactions to consider too:

- *Sensors*: The typical sensor is a camera or a Lidar (light detection and ranging), which uses a laser scanner to create 3D images. But robots might also have systems for sound, touch, taste, and even smell. In fact, they could also include sensors that go beyond human capabilities, such as night vision or detecting chemicals. The information from the sensors is sent to a controller that can activate an arm or other parts of the robot.

- *Actuators*: These are electro-mechanical devices like motors. For the most part, they help with the movement of the arms, legs, head, and any other movable part.

- *Computer*: There are memory storage and processors to help with the inputs from the sensors. In advanced robots, there may also be AI chips or Internet connections to AI cloud platforms.

Figure 7-1 shows the interactions of these functions.



**Figure 7-1.** The general system for a physical robot

---

[3] www.nytimes.com/2019/04/03/us/robotics-wheelchair.html

There are also two main ways to operate a robot. First of all, there is remote control by a human operation. In this case, the robot is called a telerobot. Then there is the autonomous robot, which uses its own abilities to navigate—such as with AI.

So what was the first mobile, thinking robot? It was Shakey. The name was apt, as the project manager of the system, Charles Rosen, noted: "We worked for a month trying to find a good name for it, ranging from Greek names to whatnot, and then one of us said, 'Hey, it shakes like hell and moves around, let's just call it Shakey.'"[4]

The Stanford Research Institute (SRI), with funding from DARPA, worked on Shakey from 1966 to 1972. And it was quite sophisticated for the era. Shakey was large, at over five feet tall, and had wheels to move and sensors and cameras to help with touching. It was also wirelessly connected to DEC PDP-10 and PDP-15 computers. From here, a person could enter commands via teletype. Although, Shakey used algorithms to navigate its environment, even closing doors.

The development of the robot was the result of a myriad of AI breakthroughs. For example, Nils Nilsson and Richard Fikes created STRIPS (Stanford Research Institute Problem Solver), which allowed for automated planning as well as the A∗ algorithm for finding the shortest path with the least amount of computer resources.[5]

By the late 1960s, as America was focused on the space program, Shakey got quite a bit of buzz. A flattering piece in *Life* declared that the robot was the "first electronic person."[6]

But unfortunately, in 1972, as the AI winter took hold, DARPA pulled the funding on Shakey. Yet the robot would still remain a key part of tech history and was inducted into the Robot Hall of Fame in 2004.[7]

# Industrial and Commercial Robots

The first real-world use of robots had to do with manufacturing industries. But these systems did take quite a while to get adoption.

---

[4] www.computerhistory.org/revolution/artificial-intelligence-robotics/13/289

[5] https://spectrum.ieee.org/view-from-the-valley/tech-history/space-age/sri-shakey-robot-honored-as-ieee-milestone

[6] www.sri.com/work/timeline-innovation/timeline.php?timeline=computing-digital#!&innovation=shakey-the-robot

[7] www.wired.com/2013/09/tech-time-warp-shakey-robot/

The story begins with George Devol, an inventor who did not finish high school. But this was not a problem. Devol had a knack for engineering and creativity, as he would go on to create some of the core systems for microwave ovens, barcodes, and automatic doors (during his life, he would obtain over 40 patents).

It was during the early 1950s that he also received a patent on a programmable robot called "Unimate." He struggled to get interest in his idea as every investor turned him down.

However, in 1957, his life would change forever when he met Joseph Engelberger at a cocktail party. Think of it like when Steve Jobs met Steve Wozniak to create the Apple computer.

Engelberger was an engineer but also a savvy businessman. He even had a love for reading science fiction, such as Isaac Asimov's stories. Because of this, Engelberger wanted the Unimate to benefit society.

Yet there was still resistance—as many people thought the idea was unrealistic and, well, science fiction—and it took a year to get funding. But once Engelberger did, he wasted little time in building the robot and was able to sell it to General Motors (GM) in 1961. Unimate was bulky (weighing 2,700 pounds) and had one 7-foot arm, but it was still quite useful and also meant that people would not have to do inherently dangerous activities. Some of its core functions included welding, spraying, and gripping—all done accurately and on a 24/7 basis.

Engelberger looked for creative ways to evangelize his robot. To this end, he appeared on Johnny Carson's *The Tonight Show* in 1966, in which Unimate putted a golf ball perfectly and even poured beer. Johnny quipped that the machine could "replace someone's job."[8]

But industrial robots did have their nagging issues. Interestingly enough, GM learned this the hard way during the 1980s. At the time, CEO Roger Smith promoted the vision of a "lights out" factory—that is, where robots could build cars in the dark!

He went on to shell out a whopping $90 billion on the program and even created a joint venture, with Fujitsu-Fanuc, called GMF Robotics. The organization would become the world's largest manufacturer of robots.

But unfortunately, the venture turned out to be a disaster. Besides aggravating unions, the robots often failed to live up to expectations. Just some of the fiascos included robots that welded doors shut or painted themselves—not the cars!

---

[8]www.theatlantic.com/technology/archive/2011/08/unimate-robot-on-johnny-carsons-tonight-show-1966/469779/

However, the situation of GMF is nothing really new—and it's not necessarily about misguided senior managers. Take a look at Tesla, which is one of the world's most innovative companies. But CEO Elon Musk still suffered major issues with robots on his factory floors. The problems got so bad that Tesla's existence was jeopardized.

In an interview on *CBS This Morning* in April 2018, Musk said he used too many robots when manufacturing the Model 3 and this actually slowed down the process.[9] He noted that he should have had more people involved.

All this points to what Hans Moravec once wrote: "It is comparatively easy to make computers exhibit adult level performance on intelligence tests or playing checkers, and difficult or impossible to give them the skills of a one-year-old when it comes to perception and mobility."[10] This is often called the Moravec paradox.

Regardless of all this, industrial robots have become a massive industry, expanding across diverse segments like consumer goods, biotechnology/healthcare, and plastics. As of 2018, there were 35,880 industrial and commercial robots shipped in North America, according to data from the Robotic Industries Association (RIA).[11] For example, the auto industry accounted for about 53%, but this has been declining.

Jeff Burnstein, president of the Association for Advancing Automation, had this to say:

> And as we've heard from our members and at shows such as Automate, these sales and shipments aren't just to large, multinational companies anymore. Small and medium-sized companies are using robots to solve real-world challenges, which is helping them be more competitive on a global scale.[12]

At the same time, the costs of manufacturing industrial robots continue to drop. Based on research from ARK, there will be a 65% reduction by 2025—with devices averaging less than $11,000 each.[13] The analysis is based on Wright's Law, which states that for every cumulative doubling in the number of units produced, there is a consistent decline in costs in percentage terms.

OK then, what about AI and robots? Where is that status of the technology? Even with the breakthroughs with deep learning, there has generally been slow progress with using AI with robots. Part of this is due to the fact that much of the research has been focused on software-based models, such as

---

[9] www.theverge.com/2018/4/13/17234296/tesla-model-3-robots-production-hell-elon-musk
[10] www.graphcore.ai/posts/is-moravecs-paradox-still-relevant-for-ai-today
[11] www.apnews.com/b399fa71204d47199fdf4c753102e6c7
[12] www.apnews.com/b399fa71204d47199fdf4c753102e6c7
[13] https://ark-invest.com/research/industrial-robot-costs

with image recognition. But another reason is that physical robots require sophisticated technologies to understand the environment—which is often noisy and distracting—in real-time. This involves enabling simultaneous localization and mapping (SLAM) in unknown environments while simultaneously tracking the robot's location. To do this effectively, there may even need to be new technologies created, such as better neural network algorithms and quantum computers.

Despite all this, there is certainly progress being made, especially with the use of reinforcement learning techniques. Consider some of the following innovations:

- *Osaro*: The company develops systems that allow robots to learn quickly. Osaro describes this as "the ability to mimic behavior that requires learned sensor fusion as well as high level planning and object manipulation. It will also enable the ability to learn from one machine to another and improve beyond a human programmer's insights." [14] For example, one of its robots was able to learn, within only five seconds, how to lift and place a chicken (the system is expected to be used in poultry factories).[15] But the technology could have many applications, such as for drones, autonomous vehicles, and IoT (Internet of Things).

- *OpenAI*: They have created the Dactyl, which is a robot hand that has human-like dexterity. This is based on sophisticated training of simulations, not real-world interactions. OpenAI calls this "domain randomization," which presents the robot many scenarios—even those that have a very low probability of happening. With Dactyl, the simulations were able to involve about 100 years of problem solving.[16] One of the surprising results was that the system learned human hand actions that were not preprogrammed—such as sliding of the finger. Dactyl also has been trained to deal with imperfect information, say when the sensors have delayed readings, or when there is a need to handle multiple objects.

---

[14] www.osaro.com/technology
[15] www.technologyreview.com/s/611424/this-is-how-the-robot-uprising-finally-begins/
[16] https://openai.com/blog/learning-dexterity/

- *MIT*: It can easily take thousands of sample data for a robot to understand its environment, such as to detect something as simple as a mug. But according to a research paper from professors at MIT, there may be a way to reduce this. They used a neural network that focused on only a few key features.[17] The research is still in the early stages, but it could prove very impactful for robots.

- *Google*: Beginning in 2013, the company went on an M&A (mergers and acquisitions) binge for robotics companies. But the results were disappointing. Despite this, it has not given up on the business. Over the past few years, Google has focused on pursuing simpler robots that are driven by AI and the company has created a new division, called Robotics at Google. For example, one of the robots can look at a bin of items and identify the one that is requested—picking it up with a three-fingered hand—about 85% of the time. A typical person, on the other hand, was able to do this at about 80%.[18]

So does all this point to complete automation? Probably not—at least for the foreseeable future. Keep in mind that a major trend is the development of cobots. These are robots that work along with people. All in all, it is turning into a much more powerful approach, as there can be leveraging of the advantages of both machines and humans.

Note that one of the major leaders in this category is Amazon.com. Back in 2012, the company shelled out $775 million for Kiva, a top industrial robot manufacturer. Since then, Amazon.com has rolled out about 100,000 systems across more than 25 fulfillment centers (because of this, the company has seen 40% improvement in inventory capacity).[19] This is how the company describes it:

> Amazon Robotics automates fulfilment center operations using various methods of robotic technology including autonomous mobile robots, sophisticated control software, language perception, power management, computer vision, depth sensing, machine learning, object recognition, and semantic understanding of commands.[20]

Within the warehouses, robots quickly move across the floor helping to locate and lift storage pods. But people are also critical as they are better able to identify and pick individual products.

---

[17] https://arxiv.org/abs/1903.06684
[18] www.nytimes.com/2019/03/26/technology/google-robotics-lab.html
[19] https://techcrunch.com/2019/03/29/built-robotics-massive-construction-excavator-drives-itself/
[20] www.amazonrobotics.com/#/vision

Yet the setup is very complicated. For example, warehouse employees wear Robotic Tech Vests so as not to be run down by robots![21] This technology makes it possible for a robot to identify a person.

But there are other issues with cobots. For example, there is the real fear that employees will ultimately be replaced by the machines. What's more, it's natural for people to feel like a proverbial cog in the wheel, which could mean lower morale. Can people really bond with robots? Probably not, especially industrial robots, which really do not have human qualities.

# Robots in the Real World

OK then, let's now take a look at some of the other interesting use cases with industrial and commercial robots.

# Use Case: Security

Both Erik Schluntz and Travis Deyle have extensive backgrounds in the robotics industry, with stints at companies like Google and SpaceX. In 2016, they wanted to start their own venture but first spent considerable time trying to find a real-world application for the technology, which involved talking to numerous companies. Schluntz and Deyle found one common theme: the need for physical security of facilities. How could robots provide protection after 5 pm—without having to spend large amounts on security guards?

This resulted in the launch of Cobalt Robotics. The timing was spot-on because of the convergence of technologies like computer vision, machine learning, and, of course, the strides in robotics.

While using traditional security technology is effective—say with cameras and sensors—they are static and not necessarily good for real-time response. But with a robot, it's possible to be much more proactive because of the mobility and the underlying intelligence.

However, people are still in the loop. Robots can then do what they are good at, such as 24/7 data processing and sensing, and people can focus on thinking critically and weighing the alternatives.

Besides its technology, Cobalt has been innovative with its business model, which it calls Robotics as a Service (RaaS). By charging a subscription, these devices are much more affordable for customers.

---

[21] www.theverge.com/2019/1/21/18191338/amazon-robot-warehouse-tech-vest-utility-belt-safety

# Use Case: Floor-Scrubbing Robots

We are likely to see some of the most interesting applications for robots in categories that are fairly mundane. Then again, these machines are really good at handling repetitive processes.

Take a look at Brain Corp, which was founded in 2009 by Dr. Eugene Izhikevich and Dr. Allen Gruber. They initially developed their technology for Qualcomm and DARPA. But Brain has since gone on to leverage machine learning and computer vision for self-driving robots. In all, the company has raised $125 million from investors like Qualcomm and SoftBank.

Brain's flagship robot is Auto-C, which efficiently scrubs floors. Because of the AI system, called BrainOS (which is connected to the cloud), the machine is able to autonomously navigate complex environments. This is done by pressing a button, and then Auto-C quickly maps the route.

In late 2018, Brain struck an agreement with Walmart to roll out 1,500 Auto-C robots across hundreds of store locations.[22] The company has also deployed robots at airports and malls.

But this is not the only robot in the works for Walmart. The company is also installing machines that can scan shelves to help with inventory management. With about 4,600 stores across the United States, robots will likely have a major impact on the retailer.[23]

# Use Case: Online Pharmacy

As a second-generation pharmacist, TJ Parker had first-hand experience with the frustrations people felt when managing their prescriptions. So he wondered: Might the solution be to create a digital pharmacy?

He was convinced that the answer was yes. But while he had a strong background in the industry, he needed a solid tech co-founder, which he found in Elliot Cohen, an MIT engineer. They would go on to create PillPack in 2013.

The focus was to reimagine the customer experience. By using an app or going to the PillPack web site, a user could easily sign up—such as to input insurance information, enter prescription needs, and schedule deliveries. When the user received the package, it would have detailed information about dose instructions and even images of each pill. Furthermore, each of the pills included labels and were presorted into containers.

---

[22] www.wsj.com/articles/walmart-is-rolling-out-the-robots-11554782460
[23] https://techcrunch.com/2019/04/10/the-startup-behind-walmarts-shelf-scanning-robots/

To make all this a reality required a sophisticated technology infrastructure, called PharmacyOS. It also was based on a network of robots, which were located in an 80,000-square-foot warehouse. Through this, the system could efficiently sort and package the prescriptions. But the facility also had licensed pharmacists to manage the process and make sure everything was in compliance.

In June 2018, Amazon.com shelled out about $1 billion for PillPack. On the news, the shares of companies like CVS and Walgreens dropped on the fears that the e-commerce giant was preparing to make a big play for the healthcare market.

## Use Case: Robot Scientists

Developing prescription drugs is enormously expensive. Based on research from the Tufts Center for the Study of Drug Development, the average comes to about $2.6 billion per approved compound.[24] In addition, it can easily take over a decade to get a new drug to market because of the onerous regulations.

But the use of sophisticated robots and deep learning could help. To see how, look at what researchers at the Universities of Aberystwyth and Cambridge have done. In 2009, they launched Adam, which was essentially a robot scientist that helped with the drug discovery process. Then a few years later, they launched Eve, which was the next-generation robot.

The system can come up with hypotheses and test them as well as run experiments. But the process is not just about brute-force calculations (the system can screen more than 10,000 compounds per day).[25] With deep learning, Eve is able to use intelligence to better identify those compounds with the most potential. For example, it was able to show that triclosan—a common element found in toothpaste to prevent the buildup of plaque—could be effective against parasite growth in malaria. This is especially important since the disease has been becoming more resistant to existing therapies.

## Humanoid and Consumer Robots

The popular cartoon, *The Jetsons*, came out in the early 1960s and had a great cast of characters. One was Rosie, which was a robot maid that always had a vacuum cleaner in hand.

---

[24] www.policymed.com/2014/12/a-tough-road-cost-to-develop-one-new-drug-is-26-billion-approval-rate-for-drugs-entering-clinical-de.html
[25] www.cam.ac.uk/research/news/artificially-intelligent-robot-scientist-eve-could-boost-search-for-new-drugs

Who wouldn't want something like this? I would. But don't expect something like Rosie coming to a home anytime soon. When it comes to consumer robots, we are still in the early days. In other words, we are instead seeing robots that have only some human features.

Here are notable examples:

- *Sophia*: Developed by the Hong Kong–based company Hanson Robotics, this is perhaps the most famous. In fact, in late 2017 Saudi Arabia granted her citizenship! Sophia, which has the likeness of Audrey Hepburn, can walk and talk. But there are also subtleties with her actions, such as sustaining eye contact.

- *Atlas*: The developer is Boston Dynamics, which launched this in the summer of 2013. No doubt, Atlas has gotten much better over the years. It can, for example, perform backflips and pick itself up when it falls down.

- *Pepper*: This is a humanoid robot, created by SoftBank Robotics, that is focused on providing customer service, such as at retail locations. The machine can use gestures— to help improve communication—and can also speak multiple languages.

As humanoid technologies get more realistic and advanced, there will inevitably be changes in society. Social norms about love and friendship will evolve. After all, as seen with the pervasiveness with smartphones, we are already seeing how technology can change the way we relate to people, say with texting and engaging in social media. According to a survey of Millennials from Tappable, close to 10% would rather sacrifice their pinky finger than forgo their smartphone![26]

As for robots, we may see something similar. It's about social robots. Such a machine—which is life-like with realistic features and AI—could ultimately become like, well, a friend or…even a lover.

Granted, this is likely far in the future. But as of now, there are certainly some interesting innovations with social robots. One example is ElliQ, which involves a tablet and a small robot head. For the most part, it is for those who live alone, such as the elderly. ElliQ can talk but also provide invaluable assistance like give reminders for taking medicine. The system can allow for video chats with family members as well.[27]

---

[26] www.mediapost.com/publications/article/322677/one-in-10-millennials-would-rather-lose-a-finger-t.html
[27] www.wsj.com/articles/on-demand-grandkids-and-robot-pals-technology-strives-to-cure-senior-loneliness-11550898010?mod=hp_lead_pos9

Yet there are certainly downsides to social robots. Just look at the awful situation of Jibo. The company, which had raised $72.7 million in venture funding, created the first social robot for the home. But there were many problems, such as product delays and the onslaught of knock-offs. Because of all this, Jibo filed for bankruptcy in 2018, and by April the following year, the servers were shut down.[28]

Needless to say, there were many disheartened owners of Jibo, evidenced by the many posts on Reddit.

# The Three Laws of Robotics

Isaac Asimov, a prolific writer of many diverse subjects like science fiction, history, chemistry, and Shakespeare, would also have a major impact on robots. In a short story he wrote in 1942 ("Runaround"), he set forth his Three Laws of Robotics:

1.  A robot may not injure a human being or, through inaction, allow a human being to come to harm.

2.  A robot must obey the orders given to it by human beings, except where such orders would conflict with the First Law.

3.  A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

---

■ **Note**   Asimov would later add another one, the zeroth law, which stated: "A robot may not harm humanity, or, by inaction, allow humanity to come to harm." He considered this law to be the most important.

---

Asimov would write more short stories that reflected how the laws would play out in complex situations, and they would be collected in a book called *I, Robot*. All these took place in the world of the 21st century.

The Three Laws represented Asimov's reaction to how science fiction portrayed robots as malevolent. But he thought this was unrealistic. Asimov had the foresight that there would emerge ethical rules to control the power of robots.

As of now, Asimov's vision is starting to become more real—in other words, it is a good idea to explore ethical principles. Granted, this may not necessarily mean that his approach is the right way. But it is a good start, especially as robots get smarter and more personal because of the power of AI.

---

[28] https://techcrunch.com/2019/03/04/the-lonely-death-of-jibo-the-social-robot/

# Cybersecurity and Robots

Cybersecurity has not been much of a problem with robots. But unfortunately, this will not likely be the case for long. The main reason is that it is becoming much more common for robots to be connected to the cloud. The same goes for other systems, such as the Internet of Things or IoT, and autonomous cars. For example, many of these systems are updated wirelessly, which exposes them to malware, viruses, and even ransoms. Furthermore, when it comes to electric vehicles, there is also a vulnerability to attacks from the charging network.

In fact, your data could linger within a vehicle! So if it is wrecked or you sell it, the information—say video, navigation details, and contacts from paired smartphone connections—may become available to other people. A white hat hacker, called GreenTheOnly, has been able to extract this data from a variety of Tesla models at junkyards, according to CNBC.com.[29] But it's important to note that the company does provide options to wipe the data and you can opt out of data collection (but this means not having certain advantages, like over-the-air (OTA) updates).

Now if there is a cybersecurity breach with a robot, the implications can certainly be devastating. Just imagine if a hacker infiltrated a manufacturing line or a supply chain or even a robotic surgery system. Lives could be in jeopardy.

Regardless, there has not been much investment in cybersecurity for robots. So far, there are just a handful of companies, like Karamba Security and Cybereason, that are focused on this. But as the problems get worse, there will inevitably be a ramping of investments from VCs and new initiatives from legacy cybersecurity firms.

# Programming Robots for AI

It is getting easier to create intelligent robots, as systems get cheaper and there are new software platforms emerging. A big part of this has been due to the Robot Operating System (ROS), which is becoming a standard in the industry. The origins go back to 2007 when the platform began as an open source project at the Stanford Artificial Intelligence Laboratory.

Despite its name, ROS is really not a true operating system. Instead, it is middleware that helps to manage many of the critical parts of a robot: planning, simulations, mapping, localization, perception, and prototypes. ROS is also modular, as you can easily pick and choose the functions you need. The result is that the system can easily cut down on development time.

---

[29] www.cnbc.com/2019/03/29/tesla-model-3-keeps-data-like-crash-videos-location-phone-contacts.html

Another advantage: ROS has a global community of users. Consider that there are over 3,000 packages for the platform.[30]

As a testament to the prowess of ROS, Microsoft announced in late 2018 that it would release a version for the Windows operating system. According to the blog post from Lou Amadio, the principal software engineer of Windows IoT, "As robots have advanced, so have the development tools. We see robotics with artificial intelligence as universally accessible technology to augment human abilities."[31]

The upshot is that ROS can be used with Visual Studio and there will be connections to the Azure cloud, which includes AI Tools.

OK then, when it comes to developing intelligent robots, there is often a different process than with the typical approach with software-based AI. That is, there not only needs to be a physical device but also a way to test it. Often this is done by using a simulation. Some developers will even start with creating cardboard models, which can be a great way to get a sense of the physical requirements.

But of course, there are also useful virtual simulators, such as MuJoCo, Gazebo, MORSE, and V-REP. These systems use sophisticated 3D graphics to deal with movements and the physics of the real world.

Then how do you create the AI models for robots? Actually, it is little different from the approach with software-based algorithms (as we covered in Chapter 2). But with a robot, there is the advantage that it will continue to collect data from its sensors, which can help evolve the AI.

The cloud is also becoming a critical factor in the development of intelligent robots, as seen with Amazon.com. The company has leveraged its hugely popular AWS platform with a new offering, called AWS RoboMaker. By using this, you can build, test, and deploy robots without much configuration. AWS RoboMaker operates on ROS and also allows the use of services for machine learning, analytics, and monitoring. There are even prebuilt virtual 3D worlds for retail stores, indoor rooms, and race tracks! Then once you are finished with the robot, you can use AWS to develop an over-the-air (OTA) system for secure deployment and periodic updates.

And as should be no surprise, Google is planning on releasing its own robot cloud platform (it's expected to launch in 2019).[32]

---

[30] www.ros.org/is-ros-for-me/
[31] https://blogs.windows.com/windowsexperience/2018/09/28/bringing-the-power-of-windows-10-to-the-robot-operating-system/
[32] www.therobotreport.com/google-cloud-robotics-platform/

# The Future of Robots

Rodney Brooks is one of the giants of the robotics industry. In 1990, he co-founded iRobot to find ways to commercialize the technology. But it was not easy. It was not until 2002 that the company launched its Roomba vacuuming robot, which was a big hit with consumers. As of this writing, iRobot has a market value of $3.2 billion and posted more than $1 billion in revenues for 2018.

But iRobot was not the only startup for Brooks. He would also help to launch Rethink Robotics—and his vision was ambitious. Here's how he put it during 2010, when his company announced a $20 million funding:

> Our robots will be intuitive to use, intelligent and highly flexible. They'll be easy to buy, train, and deploy and will be unbelievably inexpensive. [Rethink Robotics] will change the definition of how and where robots can be used, dramatically expanding the robot marketplace.[33]

But unfortunately, as with iRobot, there were many challenges. Even though Brook's idea for cobots was pioneering—and would ultimately prove to be a lucrative market—he had to struggle with the complications of building an effective system. The focus on safety meant that precision and accuracy was not up to the standards of industrial customers. Because of this, the demand for Rethink's robots was tepid.

By October 2018, the company ran out of cash and had to close its doors. In all, Rethink had raised close to $150 million from VCs and strategic investors like Goldman Sachs, Sigma Partners, GE, and Bezos Expeditions. The company's intellectual property was sold off to a German automation firm, HAHN Group.

True, this is just one example. But then again, it does show that even the smartest tech people can get things wrong. And more importantly, the robotics market has unique complexities. When it comes to the evolution of this category, progress may be choppy and volatile.

As Cobalt's Schluntz has noted:

> While the industry has made progress in the last decade, robotics hasn't yet realized its full potential. Any new technology will create a wave of numerous new companies, but only a few will survive and turn into lasting businesses. The Dot-Com bust killed the majority of internet companies, but Google, Amazon, and Netflix all survived. What robotics companies need to do is to be upfront about what their robots can do for customers today, overcome Hollywood stereotypes of robots as the bad guys, and demonstrate a clear ROI (Return On Investment) to customers.[34]

---

[33] www.rethinkrobotics.com/news-item/heartland-robotics-raises-20-million-in-series-b-financing/
[34] From the author's interview with Erik Schluntz, CTO of Cobalt Robotics.

# Conclusion

Until the past few years, robots were mostly for high-end manufacturing, such as for autos. But with the growth in AI and the lower costs for building devices, robots are becoming more widespread across a range of industries. As seen in this chapter, there are interesting use cases with robots that do things like clean floors or provide security for facilities.

But the use of AI with robotics is still in the nascent stages. Programming hardware systems is far from easy, and there is the need of sophisticated systems to navigate environments. However, with AI approaches like reinforcement learning, there has been accelerated progress.

But when thinking of using robots, it's important to understand the limitations. There also must be a clear-cut purpose. If not, a deployment can easily lead to a costly failure. Even some of the world's most innovative companies, like Google and Tesla, have had challenges in working with robots.

# Key Takeaways

- A robot can take actions, sense its environment, and have some level of intelligence. There are also key functions like sensors, actuators (such as motors), and computers.

- There are two main ways to operate a robot: the telerobot (this is controlled by a human) and autonomous robot (based on AI systems).

- Developing robots is incredibly complicated. Even some of the world's best technologists, like Tesla's Elon Musk, have had major troubles with the technology. A key reason is the Moravec paradox. Basically, what's easy for humans is often difficult for robots and vice versa.

- While AI is making an impact on robots, the process has been slow. One reason is that there has been more emphasis on software-based technologies. But also robots are extremely complicated when it comes to moving and understanding the environment.

- Cobots are machines that work alongside humans. The idea is that this will allow for the leveraging of the advantages of both machines and people.

- The costs of robots are a major reason for lack of adoption. But innovative companies, like Cobalt Robotics, are using new business models to help out, such as with subscriptions.

- Consumer robots are still in the initial stages, especially compared to industrial robots. But there are some interesting use cases, such as with machines that can be companions for people.

- During the 1950s, science fiction writer Isaac Asimov created the Three Laws of robotics. For the most part, they focused on making sure that the machines would not harm people or society. Even though there are criticisms of Asimov's approach, they are still widely accepted.

- Security has not generally been a problem with robots. But this will likely change—and fast. After all, more robots are connected to the cloud, which allows for the intrusion of viruses and malware.

- The Robot Operating System (ROS) has become a standard for the robotics industry. This middleware helps with planning, simulations, mapping, localization, perception, and prototypes.

- Developing intelligent robots has many challenges because of the need to create physical systems. Although, there are tools to help out, such as by allowing for sophisticated simulations.

# Implementation of AI

## Moving the Needle for Your Company

In March 2019, a shooter live-streamed on Facebook his brutal killing of 50 people in two mosques in New Zealand. It was viewed about 4,000 times and was not shut off until 29 minutes after the attack.[1] The video was then uploaded to other platforms and was viewed millions of times.

Yes, this was a stark example of how AI can fail in a horrible way.

In a blog post, Facebook's VP of Product Management, Guy Rosen, noted:

> AI systems are based on 'training data,' which means you need many thousands of examples of content in order to train a system that can detect certain types of text, imagery or video. This approach has worked very well for areas such as nudity, terrorist propaganda and also graphic violence where there is a large number of examples we can use to train our systems. However, this particular video did not trigger our automatic detection systems. To achieve that we will need to provide our systems with large volumes of

---

[1] www.cnbc.com/2019/03/21/why-facebooks-ai-didnt-detect-the-new-zealand-mosque-shooting-video.html

> data of this specific kind of content, something which is difficult as these events are thankfully rare. Another challenge is to automatically discern this content from visually similar, innocuous content—for example if thousands of videos from live-streamed video games are flagged by our systems, our reviewers could miss the important real-world videos where we could alert first responders to get help on the ground.[2]

It also did not help that there were various bad actors that re-uploaded edited versions of the video in order to foil Facebook's AI system.

Of course, this was a big lesson in the shortfalls of technology, and the company says it is committed to keep improving its systems. But the Facebook case study also highlights that even the most technologically sophisticated companies have major challenges. This is why when it comes to implementing AI, there needs to be solid planning as well as an understanding that there will inevitably be problems. But it can be tough as senior managers at companies are under pressure to get results from this technology.

In this chapter, we'll take a look at some of the best practices for AI implementations.

# Approaches to Implementing AI

Using AI in a company generally involves two approaches: using vendor software or creating in-house models. The first one is the most prevalent—and may be enough for a large number of companies. The irony is that you may already be using software, say from Salesforce.com, Microsoft, Google, Workday, Adobe, or SAP, that already has powerful AI capabilities. In other words, a good approach is to make sure you are taking advantage of these to the fullest.

To see what's available, take a look at Salesforce.com's Einstein, which was launched in September 2016. This AI system is seamlessly embedded into the main CRM (Customer Relationship Management) platform, allowing for more predictive and personalized actions for sales, service, marketing, and commerce. Salesforce.com calls Einstein a "personal data scientist" as it is fairly easy to use, such as with drag and drop to create the workflows. Some of the capabilities include the following:

- *Predictive Scoring*: This shows the likelihood that a lead will convert into an opportunity.

---

[2] https://newsroom.fb.com/news/2019/03/technical-update-on-new-zealand/

- *Sentiment Analysis*: This provides a way to get a sense of how people view your brand and products by analyzing social media.

- *Smart Recommendations*: Einstein crunches data to show what products are the most ideal for leads.

However, while these prebuilt features make it easier to use AI, there are still potential issues. "We have been building AI functions into our applications during the past few years and this has been a great learning experience," said Ricky Thakrar, who is Zoho's customer experience evangelist. "But to make the technology work, the users must use the software right. If the sales people are not inputting information correctly, then the results will likely be off. We also found that there should be at least three months of usage for the models to get trained. And besides, even if your employees are doing everything right, this does not mean that the AI predictions will be perfect. Always take things with a grain of salt."[3]

Now as for building your own AI models, this is a significant commitment for a company. And this is what we'll be covering in this chapter.

But regardless of what approach you may take, the implementation and use of AI should first begin with education and training. It does not matter whether the employees are non-technical people or software engineers. For AI to be successful in an organization, everyone must have a core understanding of the technology. Yes, this book will be helpful but there are many online resources to help out as well, such as from training platforms like Lynda, Udacity, and Udemy. They provide hundreds of high-quality courses on many topics about AI.

To give a sense of what a corporate training program looks like, consider Adobe. Even though the company has incredibly talented engineers, there are still a large number who do not have a background in AI. Some of them may not have specialized in this in school or their work. Yet Adobe wanted to ensure that all the engineers had a solid grasp of the core principles of AI. To this end, the company has a six-month certification program, which trained 5,000 engineers in 2018. The goal is to unleash the data scientist in each engineer.

The program includes both online courses and in-person sessions, which not only cover technical topics but also areas like strategy and even ethics. Adobe also provides help from senior computer scientists to assist students to master the topics.

Next, early on in the implementation process, it's essential to think about the potential risks. Perhaps one of the most threatening is bias since it can easily seep into an AI model.

---

[3] This is based on the author's interview, in April 2019, with Ricky Thakrar, who is Zoho's customer experience evangelist.

An example of this is Amazon.com, which shut down its AI-powered recruiting software in 2017. The main issue was that it was biased for hiring males. Interestingly enough, this was a classic case of a training problem for the model. Consider that a majority of the resume submissions were from men—so the data was skewed. Amazon.com even tried to tweak the model, but still the results were far from being gender neutral.[4]

In this case, the issue was not just about making decisions that were based on faulty premises. Amazon.com was also probably exposing itself to potential legal liability, such as with discrimination claims.

Given the tricky issues with AI, more companies are putting together ethics boards. But even this can be fraught with problems. Hey, what may be ethical for one person may not be a big deal for someone else, right? Definitely.

For example, Google closed down its own ethics board in about a week of its launch. It appears the main reason was the backlash that came from including a member from the Heritage Foundation, which is a conservative think tank.[5]

# The Steps for AI Implementation

If you plan to implement your own AI models, what are the main steps to consider? What are the best practices? Well, first of all, it's critically important that your data is fairly clean and structured in a way to allow for modelling (see Chapter 2).

Here are some other steps to look at:

- Identify a problem to solve.
- Put together a strong team.
- Select the right tools and platforms.
- Create the AI model (we went through this process in Chapter 3).
- Deploy and monitor the AI model.

Let's take a look at each.

---

[4] www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scraps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G
[5] www.theverge.com/2019/4/4/18296113/google-ai-ethics-board-ends-controversy-kay-coles-james-heritage-foundation

# Identify a Problem to Solve

Founded in 1976, HCL Technologies is one of the largest IT consulting firms, with 132,000 employees across 44 countries, and has half the Fortune 500 as customers. The company also has implemented a large number of AI systems.

Here's what Kalyan Kumar, who is the corporate vice president and global CTO of HCL Technologies, has to say:

> Business leaders need to understand and realize that the adoption of Artificial Intelligence is a journey and not a sprint. It is critical that the people driving AI adoption within an enterprise remain realistic about the timeframe and what AI is capable of doing. The relationship between humans and AI is mutually empowering, and any AI implementation may take some time before it starts to make a positive and significant impact.[6]

It's great advice. This is why—especially for companies that are starting in the AI journey—it's essential to take an experimental approach. Think of it as putting together a pilot program—that is, you are in the "crawl and walk phase."

But when it comes to the AI implementation process, it's common to get too focused on the different technologies, which are certainly fascinating and powerful. Yet success is far more than just technology; in other words, there must first be a clear business case. So here are some areas to think about when starting out:

- No doubt, decisions in companies are often ad hoc and, well, a matter of guessing! But with AI, you have an opportunity to use data-driven decision-making, which should have more accuracy. Then where in your organization can this have the biggest benefit?

- As seen with Robotic Process Automation (RPA), which we covered in Chapter 5, AI can be extremely effective when handling repetitive and mundane tasks.

- Chatbots can be another way to start out with AI. They are relatively easy to set up and can serve specific use cases, such as customer service. You can learn more about this in Chapter 6.

---

[6] This is based on the author's interview, in March 2019, with Kalyan Kumar, who is the corporate vice president and global CTO of HCL Technologies.

Andrew Ng, who is the CEO of Landing AI and the former head of Google Brain, has come up with various approaches to think about when identifying what to focus on with your initial AI project:[7]

- *Quick Win*: A project should take anywhere from 6 to 12 months and must have a high probability of success, which should help provide momentum for more initiatives. Andrew suggests having a couple projects as it increases the odds of getting a win.

- *Meaningful*: A project does not have to be transformative. But it should have results that help improve the company in a notable way, creating more buy-in for additional AI investments. The value usually comes from lower costs, higher revenues, finding new extensions of the business, or mitigating risks.

- *Industry-Specific Focus*: This is critical since a successful project will be another factor in boosting buy-in. Thus, if you have a company that sells a subscription service, then an AI system to lessen churn would be a good place to start.

- *Data*: Do not limit your options based on the amount of data you have. Andrew notes that a successful AI project may have as little as 100 data points. But the data must still be high quality and fairly clean, which are key topics covered in Chapter 2.

When looking at this phase, it is also worth evaluating the "tango" between employees and machines. Keep in mind that this is often missed—and it can have adverse consequences on an AI project. As we've seen in this book, AI is great at processing huge amounts of data with little error at great speed. The technology is also excellent with predictions and detecting anomalies. But there are tasks that humans do much better, such as being creative, engaging in abstraction, and understanding concepts.

Note the following example of this from Erik Schluntz, who is the co-founder and CTO at Cobalt Robotics:

> Our security robots are excellent at detecting unusual events in workplace and campus settings, like spotting a person in a dark office with AI-powered thermal-imaging. But one of our human operators then steps

---

[7] https://hbr.org/2019/02/how-to-choose-your-first-ai-project

in and makes the call of how to respond. Even with all of AI's potential, it's still not the best mission-critical option when pitted against constantly changing environmental variables and human unpredictability. Consider the gravity of AI making a mistake in different situations—failing to detect a malicious intruder is much worse than accidentally sounding a false alarm to one of our operators.[8]

Next, make sure you are clear-cut about the KPIs and measure them diligently. For example, if you are developing a custom chatbot for customer service, you might want to measure against metrics like the resolution rate and customer satisfaction.

And finally, you will need to do an IT assessment. If you have mostly legacy systems, then it could be more difficult and expensive to implement AI, even if vendors have APIs and integrations. This means you will need to temper your expectations.

Despite all this, the investments can truly move the needle, even for old-line companies. To see an example of this, consider Symrise, whose roots go back more than 200 years in Germany. As of this writing, the company is a global producer of flavors and fragrances, with over 30,000 products.

A few years ago, Symrise embarked on a major initiative, with the help of IBM, to leverage AI to create new perfumes. The company not only had to retool its existing IT infrastructure but also had to spend considerable time fine-tuning the models. But a big help was that it already had an extensive dataset, which allowed for more precision. Note that even a slight deviation in the mixture of a compound can make a perfume fail.

According to Symrise's president of Scent and Care, Achim Daub:

Now our perfumers can work with an AI apprentice by their side, that can analyze thousands of formulas and historical data to identify patterns and predict novel combinations, helping to make them more productive, and accelerate the design process by guiding them toward formulas that have never been seen before.[9]

---

[8] This is based on the author's interview, in April 2019, with Erik Schluntz, who is the co-founder and CTO at Cobalt Robotics.
[9] www.symrise.com/newsroom/article/breaking-new-fragrance-ground-with-artificial-intelligence-ai-ibm-research-and-symrise-are-workin/

# Forming the Team

How large should the initial team be for an AI project? Perhaps a good guide is to use Jeff Bezos' "two pizza rule."[10] In other words, is this enough to feed the people who are participating?

Oh, and there should be no rush to build the team. Everyone must be highly focused on success and understand the importance of the project. If there is little to show from the AI project, the prospects for future initiatives could be in jeopardy.

The team will need a leader who generally has a business or operational background but also has some technical skills. Such a person should be able to identify the business case for the AI project but also communicate the vision to multiple stakeholders in the company, such as the IT department and senior management.

In terms of the technical people, there will probably not be a need for a PhD in AI. While such people are brilliant, they are often focused primarily on innovations in the field, such as by refining models or creating new ones. These skillsets are usually not essential for an AI pilot.

Rather, look for those people who have a background in software engineering or data science. However, as noted earlier in the chapter, these people may not have a strong background in AI. Because of this, there may be a need to have them spend a few months of training on learning the core principles of machine learning and deep learning. There should also be a focus on understanding how to use AI platforms, such as TensorFlow.

Given the challenges, it may be a good idea to seek the help of consultants, who can help identify the AI opportunities but also provide advice on data preparation and the development of the models.

Since an AI pilot will be experimental, the team should have people who are willing to take risks and are open minded. If not, progress could be extremely difficult.

# The Right Tools and Platforms

There are many tools for helping create AI models, and most of them are open source. Even though it's good to test them out, it is still advisable to first conduct your IT assessment. By doing this, you should be in a better position to evaluate the AI Tools.

---

[10] www.geekwire.com/2018/amazon-tops-600k-worldwide-employees-1st-time-13-jump-year-ago/

Something else: You may realize that your company is already using multiple AI Tools and platforms! This may cause issues with integration and the management of the process with AI projects. In light of this, a company should develop a strategy for the tools. Think of it as your AI Tools stack.

OK then, let's take a look at some of the more common languages, platforms, and tools for AI.

# Python Language

Guido van Rossum, who got his master's degree in mathematics and computer science from the University of Amsterdam in 1982, would go on to work at various research institutes in Europe like the Corporation for National Research Initiatives (CNRI). But it was in the late 1980s that he would create his own computer language, called Python. The name actually came from the popular British comedy series *Monty Python*.

So the language was kind of offbeat—but this made it so powerful. Python would soon become the standard for AI development.

Part of this was due to the simplicity. With just a few scripts of code, you can create sophisticated models, say with functions like filter, map, and reduce. But of course, the language allows for much sophisticated coding as well.

Van Rossum developed Python with a clear philosophy:[11]

- Beautiful is better than ugly.

- Explicit is better than implicit.

- Simple is better than complex.

- Complex is better than complicated.

- Flat is better than nested.

- Sparse is better than dense.

These are just some of the principles.

What's more, Python had the advantage of growing in the academic community, which had access to the Internet that helped accelerate the distribution. But it also made it possible for the emergence of a global ecosystem with thousands of different AI packages and libraries. Here are just some:

---

[11] www.python.org/dev/peps/pep-0020/

- *NumPy*: This allows for scientific computing applications. At the heart of this is the ability to create a sophisticated array of objects at high performance. This is critical for high-end data processing in AI models.

- *Matplotlib*: With this, you can plot datasets. Often Matplotlib is used in conjunction with NumPy/Pandas (Pandas refers to "Python Data Analysis Library"). This library makes it relatively easy to create data structures for developing AI models.

- *SimpleAI*: This is an implementation of the AI algorithms from the book *Artificial Intelligence: A Modern Approach*, by Stuart Russel and Peter Norvig. The library not only has rich functionality but also provides helpful resources to navigate the process.

- *PyBrain*: This is a modular machine learning library that makes it possible to create sophisticated models—neural networks and reinforcement learning systems—without much coding.

- *Scikit-Learn*: Launched in 2007, this library has a deep source of capabilities, allowing for regression, clustering, and classification of data.

Another benefit for Python is that there are many resources for learning. A quick search on YouTube will show thousands of free courses.

Now there are other solid languages you can use for AI like C++, C#, and Java. While they are generally more powerful than Python, they are also complex. Besides, when it comes to building models, there is often little need to create full-fledged applications. And finally, there are Python libraries built for high-speed AI machines—with GPUs—like CUDA Python.

## AI Frameworks

There are a myriad of AI frameworks, which provide end-to-end systems to build models, train them, and deploy them. By far the most popular is TensorFlow, which is backed by Google. The company started development of this framework in 2011, through its Google Brain division. The goal was to find a way to create neural networks faster so as to embed the technology across many Google applications

By 2015, Google decided to open source TensorFlow, primarily because the company wanted to accelerate the progress of AI. And no doubt, this is what happened. By open sourcing TensorFlow, Google made its technology an industry standard for development. The software has been downloaded over

41 million times, and there are more than 1,800 contributors. In fact, TensorFlow Lite (which is for embedded systems) is running on more than 2 billion mobile devices.[12]

The ubiquity of the platform has resulted in a large ecosystem. This means there are many add-ons like TensorFlow Federated (for decentralized data), TensorFlow Privacy, TensorFlow Probability, TensorFlow Agents (for reinforcement learning), and Mesh TensorFlow (for massive datasets).

To use TensorFlow, you have the option of a variety of languages to create your models, such as Swift, JavaScript, and R. Although, for the most part, the most common one is Python.

In terms of the basic structure, TensorFlow takes in input data as a multidimensional array, which is also known as a tensor. There is a flow to it, represented by a chart, as the data courses through the system.

When you enter commands into TensorFlow, they are processed using a sophisticated C++ kernel. This allows for much higher performance, which can be essential as some models can be massive.

TensorFlow can be used for just about anything when it comes to AI. Here are some of the models that it has powered:

- Researchers from NERSC (National Energy Research Scientific Computing Center) at the Lawrence Berkeley National Laboratory created a deep learning system to better predict extreme weather. It was the first such model that broke the expo (1 billion billion calculations) computing barrier. Because of this, the researchers won the Gordon Bell Prize.[13]

- Airbnb used TensorFlow to build a model that categorized millions of listing photos, which increased the guest experience and led to higher conversions.[14]

- Google used TensorFlow to analyze data from NASA's Kepler space telescope. The result? By training a neural network, the model discovered two exoplanets. Google also made available the code to the public.[15]

Google has been working on TensorFlow 2.0, and a key focus is to make the API process simpler. There is also something called Datasets, which helps to streamline the preparation of data for AI models.

---

[12] https://medium.com/tensorflow/recap-of-the-2019-tensorflow-dev-summit-1b5ede42da8d
[13] www.youtube.com/watch?v=p45kQklIsd4&feature=youtu.be
[14] www.youtube.com/watch?v=tPb2u9kwh2w&feature=youtu.be
[15] https://ai.googleblog.com/2018/03/open-sourcing-hunt-for-exoplanets.html

Then what are some of the other AI frameworks? Let's take a look:

- *PyTorch*: Facebook is the developer of this platform, which was released in 2016. Like TensorFlow, the main language to program the system is Python. While PyTorch is still in the early phases, it is already considered the runner-up to TensorFlow in terms of usage. So what is different with this platform? PyTorch has a more intuitive interface. The platform also allows for dynamic computation of graphs. This means you can easily make changes to your models in runtime, which helps speed up development. PyTorch also makes it possible for having different types of back-end CPUs and GPUs.

- *Keras*: Even though TensorFlow and PyTorch are for experienced AI experts, Keras is for beginners. With a small amount of code—in Python—you can create neural networks. In the documentation, it notes: "Keras is an API designed for human beings, not machines. It puts user experience front and center. Keras follows best practices for reducing cognitive load: it offers consistent and simple APIs, it minimizes the number of user actions required for common use cases, and it provides clear and actionable feedback upon user error."[16] There is a "Getting Started" guide that takes only 30 seconds! Yet the simplicity does not mean that it is not powerful. The fact is that you can create sophisticated models with Keras. For example, TensorFlow has integrated Keras on its own platform. Even for those who are pros at AI, the system can be quite useful for doing initial experimentations with models.

With AI development, there is another common tool: Jupyter Notebook. It's not a platform or development tool. Instead, Jupyter Notebook is a web app that makes it easy to code in Python and R to create visualizations and import AI systems. You can also easily share your work with other people, similar to what GitHub does.

During the past few years, there has also emerged a new category of AI Tools called automated machine learning or autoML. These systems help to deal with processes like data prep and feature selection. For the most part, the goal is to provide help for those organizations that do not have experienced data scientists and AI engineers. This is all about the fast-growing trend of the "citizen data scientist"—that is, a person who does not have a strong technical background who can still create useful models.

---

[16] https://keras.io/

Some of the players in the autoML space include H2O.ai, DataRobot, and SaaS. The systems are intuitive and use drag-and-drop ease with the development of models. As should be no surprise, mega tech operators like Facebook and Google have created autoML systems for their own teams. In the case of Facebook, it has Asimo, which helps manage the training and testing of 300,000 models every month.[17]

For a use case of autoML, take a look at Lenovo Brazil. The company was having difficulty creating machine learning models to help predict and manage the supply chain. It had two people who coded 1,500 lines of R code each week—but this was not enough. The fact is that it would not be cost-effective to hire more data scientists.

Hence the company implemented DataRobot. By automating various processes, Lenovo Brazil was able to create models with more variables, which led to better results. Within a few months, the number of users of DataRobot went from two to ten.

Table 8-1 shows some other results.[18]

**Table 8-1.** The results of implementing an autoML system

| Tasks | Before | After |
|---|---|---|
| Model creation | 4 weeks | 3 days |
| Production models | 2 days | 5 minutes |
| Accuracy of predictions | <80% | 87.5% |

Pretty good, right? Absolutely. But there are still come caveats. With Lenovo Brazil, the company had the benefit of skilled data scientists, who understood the nuances of creating models.

However, if you use an autoML tool without such expertise, you could easily run into serious trouble. There's a good chance that you may create models that have faulty assumptions or data. If anything, the results may ultimately prove far worse than not using AI! Because of this, DataRobot actually requires that a new customer have a dedicated field engineer and data scientist work with the company for the first year.[19]

Now there are also low-code platforms that have proven to be useful in accelerating the development of AI projects. One of the leaders in the space is Appian, which has the bold guarantee of "Idea to app in eight weeks."

[17] www.aimlmarketplace.com/technology/machine-learning/the-rise-of-automated-machine-learning
[18] https://3gp10c1vpy442j63me73gy3s-wpengine.netdna-ssl.com/wp-content/uploads/2018/08/Lenovo-Case-Study.pdf
[19] www.wsj.com/articles/yes-you-too-can-be-an-ai-expert-11554168513

With this platform, you can easily set up the data structure that is clean. There are even systems in place to help guide the process, such as alerting for issues. No doubt, this provides a solid foundation for building a model. But low-code also helps in other ways. For example, you can test various AI platforms—say from Google, Amazon, or Microsoft—to see which one performs better. Then you can create the app with a modern interface and deploy it to the Web or mobile apps.

To get a sense of the power of low-code, take a look at what KPMG has done with the technology. The company was able to help its clients transition away from the use of LIBOR in loans. First of all, KPMG used its own AI platform, called Ignite, to ingest the unstructured data and use machine learning and Natural Language Processing to remediate the contracts. Next, the company used Appian to help with document sharing, customizable business rules, and real-time reporting.

Such a process—when done manually—could easily take thousands of hours, with the error rate of 10% to 15%. But when using Ignite/Appian, the accuracy was over 96%. Oh, and the time to process the documents was in seconds.

# Deploy and Monitor the AI System

Even when you build an AI model that works, there is still more work to do. You need to find ways to deploy and monitor it.

This requires change management, which is always complex and difficult. AI is different than a typical IT implementation since it involves using predictions and insights for decision-making. This means people will need to rethink how they interact with the technology.

Also consider that the chances are that the end-users will be non-technical people, whether employees or consumers. This is why there needs to be much work on making the AI model as easy as possible. For example, if you have built a system for online marketing, you might want to limit the options for the user—say to just four or five of them.

Why? If there are too many, then users may get frustrated and not even know where to start. This is all part of the so-called "analysis paralysis" problem. When this happens, there will inevitably be little adoption of the AI model, which will severely result in impeded progress.

Another good strategy is to use visualizations that are interactive. In other words, you can easily see how the trends change by adjusting some variables. You can also allow for clicking a certain part of the chart to drill down into more details.

It's also essential to create documentation. But this should be more than just written materials. For example, an effective approach is to develop video tutorials. Such an effort will go a long way in creating strong adoption.

As a best practice, the initial deployment should be limited. Perhaps this could be to a small group of beta users and a small section of the customer base. There should also be warnings that the AI model is in the early stages and may have bugs.

Therefore, this phase is about learning. What works? What should be removed? Where can things be improved?

This is definitely an iterative process that must not be rushed.

Then once the AI model is ready for full deployment, there should be enough support in place and someone to lead the management of the project. There also must be recognition for the team for the win. Hopefully, the praise will come from the highest levels of the company, which will help encourage more and more innovation.

There are a variety of automated platforms to help streamline the workflow process, such as Alteryx. The company's vision is to democratize data science and analytics, regardless if someone has a technical background or not. The Alteryx system handles the key areas of the process: data discovery, data preparation, analytics, and deployment. And all of this is done with code-free drag-and-drop tools. Furthermore, many of the company's customers are non-technology operators like Hyatt, Unilever, and Kroger.

Again, AI development is really a journey—and your strategy will inevitably change. This is inevitable. According to Kurt Muehmel, who is the VP of Sales Engineering at Dataiku[20]:

> What businesses sometimes fail to realize is that the path to AI is a long-term evolution of not only technology but in the way the company collaborates and works together. So, in addition to education, one of the key components to an AI strategy should be overall change management. It is important to create both short- and long-term roadmaps of what will be accomplished with first maybe predictive analytics, then perhaps machine learning, and ultimately—as a longer-term goal—AI, and how each roadmap impacts various pieces of the business as well as people who are a part of those business lines and their day-to-day work.

[20] This is from the author's interview, in April 2019, with Kurt Muehmel, who is the VP Sales Engineering of Dataiku.

# Conclusion

As shown in this chapter, when approaching implementing AI, it's critical to look at two paths. The first is to get the maximum use of any third-party systems that use the technology. But there should also be a focus on data quality. If not, the results will likely be off the mark.

The second path is to do an AI project, which is based on your company's own data. To be successful, there must be a strong team that has a blend of technical, business, and domain expertise. There will also likely be a need for some AI training. This is the case even for those with backgrounds in data science and engineering.

From here, there should be no rush in the steps of the project: assessing the IT environment, setting up a clear business objective, cleaning the data, selecting the right tools and platforms, creating the AI model, and deploying the system. With early projects, there will inevitably be challenges so it's critical to be flexible. But the effort should be well worth it.

# Key Takeaways

- Even the best companies have difficulties with implementing AI. Because of this, there must be great care, diligence, and planning. It's also important to realize that failure is common.

- There are two main ways to use AI in a company: through a vendor's software application or an in-house model. The latter is much more difficult and requires a major commitment from the organization.

- When using off-the-shelf AI applications, there is still much work to be done. For example, if the employees are not correctly inputting the data, then the results will likely be off.

- Education is critical with an AI implementation, even for experienced engineers. There are excellent online training resources to help out with this.

- Be mindful of the risks of AI implementations, such as bias, security, and privacy.

- Some of the key parts of the AI implementation process include the following: identify a problem to solve; put together a strong team; select the right tools and platforms; create the AI model; and deploy and monitor the AI model.

- When developing a model, look at how the technology relates to people. The fact is that people can be much better at certain tasks.

- Forming the team is not easy, so do not rush the process. Have a leader who has a good business or operational background, with a mix of technical skills.

- It's good to experiment with the various AI Tools. However, before doing this, make sure you do an IT assessment.

- Some of the popular AI Tools include TensorFlow, PyTorch, Python, Keras, and the Jupyter Notebook.

- Automated machine learning or autoML tools help to deal with processes like data prep and feature selection for AI models. The focus is on those who do not have technical skills.

- Deployment of the AI model is more than just scaling. It's also critical to have the system easy to use, so as to allow for much more adoption.

# The Future of AI

## The Pros and Cons

At the Web Summit conference in late 2017, the legendary physicist Stephen Hawking offered his opinion about the future of AI. On the one hand, he was hopeful that the technology could outpace human intelligence. This would likely mean that many horrible diseases will be cured and perhaps there will be ways to deal with environmental problems, including climate change.

But there was the dark side as well. Hawking talked about how the technology had the potential to be the "worst event in the history of our civilization."[1] Just some of the problems include mass unemployment and even killer robots. Because of this, he urged for ways to control AI.

Hawking's ideas are certainly not on the fringe. Prominent tech entrepreneurs like Elon Musk and Bill Gates also have expressed deep worries about AI.

Yet there are many who are decidedly optimistic, if not exuberant. Masayoshi Son, who is the CEO of SoftBank and the manager of the $100 billion Vision venture fund, is one of them. In an interview with CNBC, he proclaimed that within 30 years, we'll have flying cars, people will be living much longer, and we'll have cured many diseases.[2] He also noted that the main focus of his fund is on AI.

---

[1] www.cnbc.com/2017/11/06/stephen-hawking-ai-could-be-worst-event-in-civilization.html
[2] www.cnbc.com/2019/03/08/softbank-ceo-ai-will-completely-change-the-way-humans-live-within-30-years.html

OK then, who is right? Will the future be dystopian or utopian? Or will it be somewhere in the middle? Well, predicting new technologies is exceedingly difficult, almost impossible. Here are some examples of forecasts that have been wide off the mark:

- Thomas Edison declared that AC (alternating current) would fail.[3]

- In his book *The Road Ahead* (published in late 1995), Bill Gates did not mention the Internet.

- In 2007, Jim Balsillie, the co-CEO of Research in Motion (the creator of the BlackBerry device), said that the iPhone would get little traction.[4]

- In the iconic science fiction movie *Blade Runner*—released in 1982 and was set in 2019—there were many predictions that were wrong like phone booths with video phones and androids (or "replicants") that were nearly indistinguishable from humans.

Despite all this, there is one thing that is certain: In the coming years, we'll see lots of innovation and change from AI. This seems inevitable, especially since there continues to be huge amounts invested in the industry.

So then, let's take a look at some of the areas that are likely to have an outsized impact on society.

## Autonomous Cars

When it comes to AI, one of the most far-reaching areas is autonomous cars. Interestingly enough, this category is not really new. Yes, it's been a hallmark of lots of science fiction stories for many decades! But for some time, there have been many real-life examples of innovation, like the following:

- *Stanford Cart*: Its development started in the early 1960s, and the original goal was to create a remote-controlled vehicle for moon missions. But the researchers eventually changed their focus and developed a basic autonomous vehicle, which used cameras and AI for navigation. While it was a standout achievement for the era, it was not practical as it required more than 10 minutes to plan for any move!

---

[3] www.msn.com/en-us/news/technology/the-best-and-worst-technology-predictions-of-all-time/ss-BBIMwm3#image=5
[4] www.recode.net/2017/1/9/14215942/iphone-steve-jobs-apple-ballmer-nokia-anniversary

- *Ernst Dickmanns*: A brilliant German aerospace engineer, he would turn his attention to the idea of converting a Mercedes van into an autonomous vehicle…in the mid-1980s. He wired together cameras, sensors, and computers. He also was creative in how he used software, such as by only focusing the graphics processing on important visual details to save on power. By doing all this, he was able to develop a system that would control a car's steering, gas pedal, and brakes. He tested the Mercedes on a Paris highway—in 1994—and it went over 600 miles, with a speed up to 81 MPH.[5] Nevertheless, the research funding was pulled because it was far from clear if there could be commercialization in a timely manner. It also did not help that AI was entering another winter.

But the inflection point for autonomous cars came in 2004. The main catalyst was the Iraq War, which was taking a horrible toll on American soldiers. For DARPA, the belief was that autonomous vehicles could be a solution.

But the agency faced many tough challenges. This is why it set up a contest, dubbed the DARPA Grand Challenge, in 2004, which had a $1 million grand prize to encourage wider innovation. The event involved a 150-mile race in the Mojave Desert, and unfortunately, it was not encouraging as the cars performed miserably. None of them finished the race!

But this only spurred even more innovation. By the next year, five cars finished the race. Then in 2007, the cars were so advanced that they were able to take actions like U-turns and merging.

Through this process, DARPA was able to allow for the creation of the key components for autonomous vehicles:

- *Sensors*: These include radar and ultrasonic systems that can detect vehicles and other obstacles, such as curbs.

- *Video Cameras*: These can detect road signs, traffic lights, and pedestrians.

- *Lidar (Light Detection and Ranging)*: This device—which is usually at the top of an autonomous car—shoots laser beams to measure the surroundings. The data is then integrated into existing maps.

- *Computer*: This helps with the control of the car, including the steering, acceleration, and braking. The system leverages AI to learn but also has built-in rules for avoiding objects, obeying the laws, and so on.

---

[5] www.politico.eu/article/delf-driving-car-born-1986-ernst-dickmanns-mercedes/

Now when it comes to autonomous cars, there is lots of confusion of what "autonomous" really means. Is it when a car drives itself completely alone—or must there be a human driver?

To understand the nuances, there are five levels of autonomy:

- *Level 0*: This is where a human controls all the systems.

- *Level 1*: With this, computers control limited functions like cruise control or braking—but only one at a time.

- *Level 2*: This type of car can automate two functions.

- *Level 3*: This is where a car automates all the safety functions. But the driver can intervene if something goes wrong.

- *Level 4*: The car can generally drive itself. But there are cases in which a human must participate.

- *Level 5*: This is the Holy Grail, in which the car is completely autonomous.

The auto industry is one of the biggest markets, and AI is likely to unleash wrenching changes. Consider that transportation is the second largest household expenditure, behind housing, and twice as large as healthcare. Something else to keep in mind: The typical car is used only about 5% of the time as it is usually parked somewhere.[6]

In light of the enormous opportunity for improvement, it should be no surprise that the autonomous car industry has seen massive amounts of investment. This has not only been about venture capitalists investing in a myriad of startups but also innovation from traditional automakers like Ford, GM, and BMW.

Then when might we see this industry become mainstream? The estimates vary widely. But according to a study from Allied Market Research, the market is forecasted to hit $556.67 billion by 2026, which would represent a compound annual growth rate of 39.47%.[7]

But there is still much to work out. "At best, we are still years away from a car that doesn't require a steering wheel," said Scott Painter, who is the CEO and founder of Fair. "Cars will still need to be insured, repaired, and maintained, even if you came back from the future in a Delorean and brought the manual for how to make these cars fully autonomous. We make 100 million cars-per-year, of which 16 million-a-year are in the U.S. And, supposing you wanted the whole supply to have these artificial intelligence features, it would still take 20 years until we had more cars on the road including all the different levels of A.I. versus the number of cars that didn't have those technologies."[8]

---

[6] www.sec.gov/Archives/edgar/data/1759509/000119312519077391/d633517ds1a.htm
[7] www.alliedmarketresearch.com/autonomous-vehicle-market
[8] From the author's interview, in May 2019, with Scott Painter, who is the CEO and founder of Fair.

But there are many other factors to keep in mind. After all, the fact remains that driving is complex, especially in urban and suburban areas. What if a traffic sign is changed or even manipulated? How about if an autonomous car must deal with a dilemma like having to decide to crash into an oncoming car or plunging into a curb, which may have pedestrians? All these are extremely difficult.

Evening seemingly simple tasks can be tough to pull off. John Krafcik, who is the CEO of Google's Waymo, points out that parking lots are a prime example.[9] They require finding available spots, avoiding other cars and pedestrians (that can be unpredictable), and moving into the space.

But technology is just one of the challenges with autonomous vehicles. Here are some others to consider:

- *Infrastructure*: Our cities and towns are built for traditional cars. But by mixing autonomous vehicles, there will probably be many logistical issues. How does a car anticipate the actions of human drivers? Actually, there may be a need to install sensors alongside roads. Or another option is to have separate roads for autonomous vehicles. Governments also will probably need to change driver's ed, providing guidance on how to interact with autonomous vehicles while on the road.

- *Regulation*: This is a big wild card. For the most part, this may be the biggest impediment as governments tend to work slowly and are resistant to change. The United States is also a highly litigious country—which may be another factor that could curb development.

- *Adoption*: Autonomous vehicles will probably not be cheap, as systems like Lidar are costly. This will certainly be a limiting factor. But at the same time, there are indications of skepticism from the general public. According to a survey from AAA, about 71% of the respondents said they are afraid of riding in an autonomous vehicle.[10]

Given all this, the initial phase of autonomous vehicles will probably be for controlled situations, say for trucking, mining, or shuttles. A case of this is Suncor Energy, which uses autonomous trucks for excavating various sites in Canada.

Ride-sharing networks—like Uber and Lyft—may be another starting point. These services are fairly structured and understandable to the public.

---

[9] www.businessinsider.com/waymo-ceo-john-krafcik-explains-big-challenge-for-self-driving-cars-2019-4
[10] https://newsroom.aaa.com/2019/03/americans-fear-self-driving-cars-survey/

Keep in mind that Waymo has been testing a self-driving taxi service in Phoenix (this is similar to a ride-sharing system like Uber, but the cars have autonomous systems). Here's how a blog post from the company explains it:

> We'll start by giving riders access to our app. They can use it to call our self-driving vehicles 24 hours a day, 7 days a week. They can ride across several cities in the Metro Phoenix area, including Chandler, Tempe, Mesa, and Gilbert. Whether it's for a fun night out or just to get a break from driving, our riders get the same clean vehicles every time and our Waymo driver with over 10 million miles of experience on public roads. Riders will see price estimates before they accept the trip based on factors like the time and distance to their destination.[11]

Waymo has found that a key is education because the riders have lots of questions. To deal with this, the company has built in a chat system in the app to contact a support person. The dashboard of the car also has a screen that provides details of the ride.

According to the blog post, "Feedback from riders will continue to be vital every step of the way."[12]

# US vs. China

The rapid ascent of China has been astonishing. Within a few years, the economy may be larger than the United States, and a key part of the growth will be AI. The Chinese government has set forth the ambitious goal of spending $150 billion on this technology through 2030.[13] In the meantime, there will continue to be major investments from companies like Baidu, Alibaba, and Tencent.

Even though China is often considered to not be as creative or innovative as Silicon Valley—often tagged as "copycats"—this perception may prove to be a myth. A study from the Allen Institute for Artificial Intelligence highlights that China is expected to outrank the United States in the most cited technical papers on AI.[14]

The country has some other advantages, which AI expert and venture capitalist Kai-Fu Lee has pointed out in his provocative book, *AI Superpowers: China, Silicon Valley, and the New World Order*[15]:

---

[11] https://medium.com/waymo/riding-with-waymo-one-today-9ac8164c5c0e
[12] Ibid.
[13] www.diamandis.com/blog/rise-of-ai-in-china
[14] www.theverge.com/2019/3/14/18265230/china-is-about-to-overtake-america-in-ai-research
[15] New York: Houghton Mifflin Harcourt, 2018.

- *Enthusiasm*: Back in the 1950s, Russia's launch of Sputnik sparked interest in people in the United States to become engineers for the space program. Something similar has actually happened in China. When the country's top Go player, Ke Jie, lost to the AlphaGo AI system, this was a wake-up call. The result is that this has inspired many young people to pursue a career in AI.

- *Data*: With a population of over 1.3 billion, China is rich with data (there are more than 700 million Internet users). But the country's authoritarian government is also critical as privacy is not considered particularly important, which means there is much more leeway when developing AI models. For example, in a paper published in *Nature Medicine*, the Chinese researchers had access to data on 600,000 patients to conduct a healthcare study.[16] While still in the early stages, it showed that an AI model was able to effectively diagnose childhood conditions like the flu and meningitis.

- *Infrastructure*: As a part of the Chinese government's investment plans, there has been a focus on creating next-generation cities that allow for autonomous cars and other AI systems. There has also been an aggressive rollout of 5G networks.

As for the United States, the government has been much more tentative with AI. President Trump has signed an executive order—called the "American AI Initiative"—to encourage development of the technology, but the terms are vague and it is far from clear how much money will be committed to it.

# Technological Unemployment

The concept of technological unemployment, which gained notoriety from famed economist John Maynard Keynes during the Great Depression, explains how innovations can lead to long-term job loss. However, evidence of this has been elusive. Notwithstanding the fact that automation has severely impacted industries like manufacturing, there is often a transition of the workforce as people adapt.

But could the AI revolution be different? It very well could. For example, California Governor Gavin Newsom fears that his state could see massive unemployment in areas like trucking and warehousing—and soon.[17]

---

[16] www.nature.com/articles/s41591-018-0335-9
[17] www.mercurynews.com/2019/03/18/were-not-prepared-for-the-promise-of-artificial-intelligence-experts-warn/

Here's another example: Harvest CROO Robotics has built a robot, called Harv, that can pick strawberries and other plants without causing bruises. Granted, it is still in the experimental phase, but the system is quickly improving. The expectation is that one robot will do the work of 30 people.[18] And of course, there will be no wages to pay or labor liability exposure.

But AI may mean more than replacing low-skilled jobs. There are already signs that the technology could have a major impact on white-collar professions. Let's face it, there is even more incentive to automate these jobs because they fetch higher compensation.

Just one category that could face AI job loss is the legal field, as a variety of startups are gunning for the market like Lawgood, NexLP, and RAVN ACE. The solutions are focused on automating areas such as legal research and contract review.[19] Even though the systems are far from perfect, they can certainly process much more volume than people—and can also get smarter as they are used more and more.

True, the overall job market is dynamic, and there will be new types of careers that will be created. There will also likely be AI innovations that are assistive for employees—making their job easier to do. For example, software startup Measure Square has been able to use sophisticated algorithms to convert paper-based floorplans into digitally interactive floorplans. Because of this, it has been easier to get projects started and completed on time.

However, in light of the potential transformative impact of AI, it does seem reasonable that there will be an adverse impact on a broad range of industries. Perhaps a foreshadowing of this is what happened with job losses from manufacturing in the 1960s to 1990s. According to the Pew Research Center, there has been virtually no real wage growth in the last 40 years.[20] During this period, the United States has also experienced a widening gap in wealth. Berkeley economist Gabriel Zucman estimates that 0.1% of the population controls nearly 20% of the wealth.[21]

---

[18] www.washingtonpost.com/news/national/wp/2019/02/17/feature/inside-the-race-to-replace-farmworkers-with-robots/
[19] www.cnbc.com/2017/02/17/lawyers-could-be-replaced-by-artificial-intelligence.html
[20] www.pewresearch.org/fact-tank/2018/08/07/for-most-us-workers-real-wages-have-barely-budged-for-decades/
[21] http://fortune.com/2019/02/08/growing-wealth-inequality-us-study/

Yet there are actions that can be taken. First of all, governments can look to provide education and transition assistance. With the pace of change in today's world, there will need to be ongoing renewal of skills for most people. IBM CEO Ginni Rometty has noted that AI will change all jobs within the next 5–10 years. By the way, her company has seen a 30% reduction of headcount in the HR department because of automation.[22]

Next, there are some people who advocate basic income, which provides a minimum amount of compensation to everyone. This would certainly soften some of the inequality, but it also has drawbacks. People definitely get pride and satisfaction from their careers. So what might a person's morale be if he or she cannot find a job? It could have a profound impact.

Finally, there is even talk of some type of AI tax. This would essentially claw back the large gains from those companies that benefit from the technology. Although, given their power, it probably would be tough to pass this type of legislation.

# The Weaponization of AI

The Air Force Research Lab is working on prototypes for something called Skyborg. It's right out of *Star Wars*. Think of Skyborg as R2-D2 that serves as an AI wingman for a fighter jet, helping to identify targets and threats.[23] The AI robot may also be able to take control if the pilot is incapacitated or distracted. The Air Force is even looking at using the technology to operate drones.

Cool, huh? Certainly. But there is a major issue: By using AI, might humans ultimately be taken out of the loop when making life-and-death decisions on the battlefield? Could this ultimately lead to more bloodshed? Perhaps the machines will make the wrong decisions—causing even more problems?

Many AI researchers and entrepreneurs are concerned. To this end, more than 2,400 have signed a statement that calls for a ban of so-called robot killers.[24]

Even the United Nations is exploring some type of ban. But the United States, along with Australia, Israel, the United Kingdom, and Russia, have resisted this move.[25] As a result, there may be a true AI arms race emerging.

---

[22] www.cnbc.com/2019/04/03/ibm-ai-can-predict-with-95-percent-accuracy-which-employees-will-quit.html
[23] www.popularmechanics.com/military/aviation/a26871027/air-force-ai-fighter-plane-skyborg/
[24] www.theguardian.com/science/2018/jul/18/thousands-of-scientists-pledge-not-to-help-build-killer-ai-robots
[25] www.theguardian.com/science/2019/mar/29/uk-us-russia-opposing-killer-robot-ban-un-ai

According to a paper from the RAND Corporation, there is even the potential that the technology could lead to nuclear war, say by the year 2040. How? The authors note that AI may make it easier to target submarines and mobile missile systems. According to the report:

> Nations may be tempted to pursue first-strike capabilities as a means of gaining bargaining leverage over their rivals even if they have no intention of carrying out an attack, researchers say. This undermines strategic stability because even if the state possessing these capabilities has no intention of using them, the adversary cannot be sure of that.[26]

But in the near term, AI will probably have the most impact on information warfare, which could still be highly destructive. We got a glimpse of this when the Russian government interfered with the 2016 presidential election. The approach was fairly low-tech as it used social media troll farms to disseminate fake news—but the consequences were significant.

But as AI gets more powerful and becomes more affordable, we'll likely see it supercharge these kinds of campaigns. For example, deepfake systems can easily create life-like photos and videos of people that could be used to quickly spread messages.

# Drug Discovery

The advances in drug discovery have been almost miraculous as we now have cures for such intractable diseases like hepatitis C and have continued to make strides with a myriad of cancers. But of course, there is certainly much that needs to be done. The fact is that drug companies are having more troubles coming up with treatments. Here's just one example: In March 2019, Biogen announced that one of its drugs for Alzheimer's, which was in Phase III trials, failed to show meaningful results. On the news, the company's shares plunged by 29%, wiping out $18 billion of market value.[27]

Consider that traditional drug development often involves much trial and error, which can be time consuming. Then might there be a better way?

Increasingly, researchers are looking to AI for help. We are seeing a variety of startups spring up that are focusing on the opportunity.

---

[26] www.rand.org/news/press/2018/04/24.html
[27] www.wsj.com/articles/biogen-shares-drop-28-after-ending-alzheimers-phase-3-trials-11553170765

One is Insitro. The company, which got its start in 2019, had little trouble raising a staggering $100 million in its Series A round. Some of the investors included Alexandria Venture Investments, Bezos Expeditions (which is the investment firm of Amazon.com's Jeff Bezos), Mubadala Investment Company, Two Sigma Ventures, and Verily.

Even though the team is relatively small—with about 30 employees—they all are brilliant researchers who span areas like data science, deep learning, software engineering, bioengineering, and chemistry. The CEO and founder, Daphne Koller, has the rare blend of experience in advanced computer science and health sciences, having led Google's healthcare business, Calico.

As a testament to Insitro's prowess, the company has already struck a partnership with mega drug operator Gilead. It involves potential payments of over $1 billion for research on nonalcoholic steatohepatitis (NASH), which is a serious liver disease.[28] A key is that Gilead has been able to assemble a large amount of data, which can train the models. This will be done using cells outside of a person's body—that is, with an in vitro system. Gilead has some urgency for looking at alternative approaches since one of its NASH treatments, selonsertib, failed in its clinical trials (it was for those who had the disease in the later stages).

The promise of AI is that it will speed up drug discovery because deep learning should be able to identify complex patterns. But the technology could also turn out to be helpful in developing personalized treatments—such as geared to a person's genetic make-up—which is likely to be critical for curing certain diseases.

Regardless, it is probably best to temper expectations. There will be major hurdles to deal with as the healthcare industry will need to undergo changes because there will be increased education for AI. This will take time, and there will likely be resistance.

Next, deep learning is generally a "black box" when it comes to understanding how the algorithms really work. This could prove difficult in getting regulatory approval for new drugs as the FDA focuses on causal relationships.

Finally, the human body is highly sophisticated, and we still are learning about how it works. And besides, as we have seen with innovations like the decoding of the Human Genome, it usually takes considerable time to understand new approaches.

As a sign of the complexities, consider the situation of IBM's Watson. Even though the company has some of the most talented AI researchers and has spent billions on the technology, it recently announced that it would no longer sell Watson for drug discovery purposes.[29]

---

[28] www.fiercebiotech.com/biotech/stealthy-insitro-opens-up-starting-gilead-deal-worth-up-to-1-05b
[29] https://khn.org/morning-breakout/ups-and-downs-of-artificial-intelligence-ibm-stops-sales-development-of-watson-for-drug-discovery-hospitals-learn-from-ehrs/

# Government

An article from Bloomberg.com in April 2019 caused a big stir. It described a behind-the-scenes look at how Amazon.com manages its Alexa speaker AI system.[30] While much of it is based on algorithms, there are also thousands of people who analyze voice clips in order to help make the results better. Often the focus is on dealing with the nuances of slang and regional dialects, which have been difficult for deep learning algorithms.

But of course, it's natural for people to wonder: Is my smart speaker really listening to me? Are my conversations private?

Amazon.com was quick to point out that it has strict rules and requirements. But even this ginned up even more concern! According to the Bloomberg.com post, the AI reviewers would sometimes hear clips that involved potentially criminal activity, such as sexual assault. But Amazon apparently has a policy to not interfere.

As AI becomes more pervasive, we'll have more of these kinds of stories; and for the most part, there will not be clear-cut answers. Some people may ultimately decide not to buy AI products. Yet this will probably be a small group. Hey, even with the myriad of privacy issues with Facebook, there has not been a decline in the user growth.

More likely, governments will start to wade in with AI issues. A group of congresspersons have sponsored a bill, called the Algorithmic Accountability Act, which aims to mandate that companies audit their AI systems (it would be for larger companies, with revenues over $50 million and more than 1 million users).[31] The law, if enacted, would be enforced by the Federal Trade Commission.

There are also legislative moves from states and cities. In 2019, New York City passed its own law to require more transparency with AI.[32] There are also efforts in Washington state, Illinois, and Massachusetts.

With all this activity, some companies are getting proactive, such as by adopting their own ethics boards. Just look at Microsoft. The company's ethics board, called Aether (AI and Ethics in Engineering and Research), decided to not allow the use of its facial recognition system for traffic stops in California.[33]

---

[30] www.bloomberg.com/news/articles/2019-04-10/is-anyone-listening-to-you-on-alexa-a-global-team-reviews-audio
[31] www.theverge.com/2019/4/10/18304960/congress-algorithmic-accountability-act-wyden-clarke-booker-bill-introduced-house-senate
[32] www.wsj.com/articles/our-software-is-biased-like-we-are-can-new-laws-change-that-11553313609?mod=hp_lead_pos8
[33] www.geekwire.com/2019/policing-ai-task-industry-government-customers/

In the meantime, we may see AI activism as well, in which people organize to protest the use of certain applications. Again, Amazon.com has been the target of this, with its Rekognition software that uses facial recognition to help law enforcement identify suspects. The ACLU has raised concerns of accuracy of the system, especially regarding women and minorities. In one of its experiments, it found that Rekognition identified 28 members of the Congress as having prior criminal records![34] As for Amazon.com, it has disputed the claims.

Rekognition is only one among various AI applications in law enforcement that are leading to controversy. Perhaps the most notable example is COMPAS (Correctional Offender Management Profiling for Alternative Sanctions), which uses analytics to gauge the probability of someone who may commit a crime. The system is often used for sentencing. But the big issue is: Might this violate a person's constitutional right to due process since there is the real risk that the AI will be incorrect or discriminatory? Actually, for now, there are few good answers. But given the importance AI algorithms will play in our justice system, it seems like a good bet that the Supreme Court will be making new law.

# AGI (Artificial General Intelligence)

In Chapter 1, we learned about the difference between strong and weak AI. And for the most part, we are in the weak AI phase, in which the technology is used for narrow categories.

As for strong AI, it's about the ultimate: the ability for a machine to rival a human. This is also known as Artificial General Intelligence or AGI. Achieving this is likely many years away, perhaps something we may not see until the next century or ever.

But of course, there are some brilliant researchers who believe that AGI will come soon. One is Ray Kurzweil, who is an inventor, futurist, bestselling author, and director of Engineering at Google. When it comes to AI, he has left his imprint on the industry, such as with innovations in areas like text-to-speech systems.

Kurzweil believes that AGI will happen—in which the Turing Test will be cracked—in 2019, and then by 2045, there will be the Singularity. This is where we'll have a world of hybrid people: part human, part machine.

Kind of crazy? Perhaps so. But Kurzweil does have many high-profile followers.

---

[34] www.businessinsider.com/ai-experts-call-on-amazon-not-to-sell-rekognition-software-to-police-2019-4

But there is much heavy lifting to be done to get to AGI. Even with the great strides with deep learning, it still generally requires large amounts of data and significant computing power.

AGI will instead need new approaches, such as the ability to use unsupervised learning. Transfer learning will likely be critical as well. For example, as we've covered earlier in the book, AI has been able to realize superhuman capabilities in playing games like Go. But transfer learning would mean that this system would be able to leverage this knowledge to play other games or to learn other fields.

In addition, AGI will need to have the capacity for common sense, abstraction, curiosity, and finding causal relationships, not just correlations. Such abilities have proven extremely difficult with computers. If anything, there will need to be breakthroughs in hardware and chip technologies. This is the opinion of Yann LeCun, one of the world's top AI researchers and the chief artificial intelligence scientist at Facebook.[35] He also thinks there needs to be much more progress with batteries and other energy sources.

Something else that will be critical: more diversity within the AI field. According to a report from the AI Now Institute, about 80% of AI professors are men; and among the AI research staffs at Facebook and Google, women accounted for 15% and 10%, respectively.[36]

This lopsidedness means that research could be more susceptible to bias. Furthermore, there will be the loss of the benefit of broader views and insights.

# Social Good

The management consulting firm, McKinsey & Co., has written an extensive study entitled "Applying Artificial Intelligence for Social Good."[37] It shows how AI is being used to deal with such issues as poverty, natural disasters, and improving education. The study has roughly 160 use cases. So here's a look at just some:

- The analysis of social media platforms can help track the outbreak of a disease.

- A nonprofit, called the Rainforest Connection, uses TensorFlow to create AI models—based on audio data— to locate illegal logging.

---

[35] http://fortune.com/2019/02/18/facebook-yann-lecun-lawnmowers-deep-learning/
[36] www.theverge.com/2019/4/16/18410501/artificial-intelligence-ai-diversity-report-facial-recognition
[37] www.mckinsey.com/featured-insights/artificial-intelligence/applying-artificial-intelligence-for-social-good

- Researchers have built a neural network that is trained on videos of poachers in Africa. With this, a drone flies over areas to detect violators, such as by using thermal infrared images.

- AI is being used to analyze data on 55,893 parcels in the city of Flint to find evidence of lead poisoning. The system relies primarily on a Bayesian model, which allows for more sophisticated predictions of toxicity. This means the health workers can more quickly take action if there are any problems in the city, potentially saving lives.

# Conclusion

This topic is a good place to end this book, I think. Regardless of all the potential for harm and the adverse consequences, AI truly has the promise for being transformative for the world. And the good news is that there are many people who are focused on making this a reality. It's not about making huge amounts of money or getting fame. The goal is to change the world.

# Key Takeaways

- Autonomous cars are far from new. But the inflection point for the development of this technology came in 2004, with a contest sponsored by DARPA.

- Some of the key components of an autonomous car include video cameras, Lidar (lasers that help process the environment), and sensors (such as for detecting other vehicles and obstacles like curbs).

- In terms of defining what is "autonomous," there are five levels. The fifth is when the vehicle is fully autonomous.

- Some of the challenges for autonomous cars are infrastructure (existing highways are not ideal), regulation, costs, and consumer adoption.

- The United States is considered the global leader in AI. But this could change soon. China is investing heavily in AI and has major advantages like enormous amounts of data and large numbers of skilled engineers.

- One of the fears of AI is that it will lead to mass unemployment, whether for blue-collar or while-collar jobs. It's true that technology has already impacted industries, like manufacturing, but the markets have proven adaptable. But if AI is transformative, it could lead to quite a bit of disruption. This is why there will likely be a need for training and re-skilling for new careers.

- Drones have had a major impact on warfare. But with AI, it's becoming possible to allow this technology to make the decisions on the battlefield. Now there are many people who see this as a big problem. However, the United States, Russia, and other countries appear to be focused on pursuing autonomous weapons.

- But when it comes to warfare—at least in the near term—AI may have a more immediate effect with the spread of false information. We saw this with the Russian's interference with the 2016 presidential election.

- AI is expected to greatly help with the drug discovery process. Already mega pharma operators, like Gilead, are exploring the technology. AI can not only process huge amounts of data but also detect patterns that may not be discernable for humans.

- As AI becomes more pervasive, there will be growing concerns about privacy and transparency. Because of this, there have been moves in the Congress, including cities and states, to impose regulations. It's not clear what may transpire, but it seems likely we'll see more restrictions. In the meantime, some companies are trying to be proactive, such as by setting up ethics boards.

- Artificial General Intelligence or AGI is where a system has human intelligence. We are likely a long way from this, though. The reason is that there will need to be new innovations in AI, such as with unsupervised learning and the creation of new hardware.

# AI Resources

## Publications and Blogs That Cover AI

- aitrends.com: www.aitrends.com/

- The Berkeley Artificial Intelligence (BAIR): https://bair.berkeley.edu/blog/

- KDnuggets: www.kdnuggets.com/news/index.html

- Machine Learning Mastery: https://machinelearning-mastery.com/blog/

- MIT Technology Review: www.technologyreview.com/

- ScienceDaily—AI Section: www.sciencedaily.com/news/computers_math/artificial_intelligence/

## Company AI Blogs

- Baidu: http://research.baidu.com/

- DeepMind: https://deepmind.com/blog/

- Facebook: https://research.fb.com/blog/

- Google: https://ai.googleblog.com/

- Microsoft: www.microsoft.com/en-us/research/

- NVIDIA: https://blogs.nvidia.com/blog/category/deep-learning/

- OpenAI: https://openai.com/blog/

# Twitter Feeds of Top AI Researchers

- Fei-Fei Li: https://twitter.com/drfeifei
- Ian Goodfellow: https://twitter.com/goodfellow_ian
- Demis Hassabis: https://twitter.com/demishassabis
- Yann Lecun: https://twitter.com/ylecun?
- Andrew Ng: https://twitter.com/AndrewYNg

# Open Source AI Tools and Platforms

- Jupyter Notebook: https://jupyter.org/
- Keras: https://keras.io/
- Python language: www.python.org/
- PyTorch: https://pytorch.org/
- TensorFlow: www.tensorflow.org/

# Online Courses

- Coursera: www.coursera.org/
- Udacity: www.udacity.com/
- Udemy: www.udemy.com/

# Glossary

*Activation Function*: Used in deep learning models to help calculate non-linear relationships.

*Actuators*: Electro-mechanical devices like motors. They help with the movement of a robot.

*AI*: See Artificial Intelligence.

*AI Winter*: A prolonged period of time, such as in the 1970s and 1980s, when the AI industry came under much pressure, such as with cutbacks in funding.

*Artificial Intelligence*: Where computers are able to learn from experience, which often involves processing data using sophisticated algorithms. Artificial intelligence is a broad category, which includes subsets like machine learning, deep learning, and Natural Language Processing (NLP).

*Artificial Neural Network (ANN)*: The most basic structure for a deep learning model. The ANN includes multiple hidden layers that process data through the use of sophisticated algorithms.

*Automation Fatigue*: With RPA, there will generally be less improvement as more tasks are automated.

*Automated Machine Learning (AutoML)*: A digital tool or platform that allows beginners to create their own AI models.

*Backpropagation*: A major breakthrough in deep learning. Backpropagation allows for more efficient assigning of weightings in models.

*Bayes' Theorem*: A statistical measure used in machine learning that helps to provide a more accurate view of the probabilities.

*Big Data*: A category of technology that involves processing huge amounts of data. Big Data is often described as having the three Vs—that is, volume, variety, and velocity.

*Binning*: Involves organizing data into groups.

*Categorical Data*: Data that does not have a numerical meaning but instead has textual meaning, say with describing race or gender.

*Cerebral Cortex*: Part of the human brain that has the most similarities to AI. It helps with thinking and other cognitive activities.

*Chatbot*: An AI system that communicates with people

*Clustering*: A form of unsupervised learning that takes unlabeled data and uses algorithms to put similar items into groups.

*Cobot*: A robot that works alongside people.

*Cognitive Robotic Process Automation (CRPA)*: An RPA system that leverages AI technologies.

*Convolutional Neural Network (CNN)*: A deep learning model that goes through different variations—or convolutions—of analysis on data. CNNs are often used for complex applications like facial recognition.

*Data Lake*: Allows for the storage and processing of massive amounts of structured and unstructured data. There is often little to no need to re-format the data.

*Data Type*: The kind of information a variable represents, such as a Boolean, integer, string, or floating point number.

*Decision Tree*: A machine learning algorithm that is a workflow of decision paths.

*Deepfake*: Involves using deep learning models to create images or videos that are misleading or harmful.

*Deep Learning*: A type of AI that uses neural networks, which mimic the processes of the brain. Much of the innovation in the field during the past decade has been with deep learning research.

*Ensemble Modelling*: Involves using more than one model for generating predictions.

*ETL (Extraction, Transformation, and Load)*: A form of data integration that is typically used in a data warehouse.

*Ethics Board*: A committee that evaluates the issues of AI projects.

*Expert System*: An early type of AI application that emerged in the 1980s. It used sophisticated logic systems to help understand certain areas like medicine, finance, and manufacturing.

*Explainability*: The process of understanding the underlying causes of a deep learning model.

*False Positive*: When a model prediction shows that the result is true even though it is not.

*Feature*: This is a column of data.

*Feature Engineering*: See Feature Extraction.

*Feature Extraction*: Describes the process of selecting the variables for an AI model.

*Feed-Forward Neural Network*: A deep learning model that processes data in a linear direction through the hidden layers. There is no cycling back.

*Generative Adversarial Network (GAN)*: Developed by AI researcher Ian Goodfellow, this is a next-generation deep learning model that helps to create new outputs like audio, text, or video.

*GPUs (Graphics Processing Units)*: Chips that were originally used for high-speed video games because of the ability to process large amounts of data quickly. But GPUs have also proven to be adept at handling AI applications.

*Hadoop*: Allows for managing Big Data, such as by making it possible to create sophisticated data warehouses.

*Hidden Layers*: The different levels of analysis in a deep learning model.

*Hidden Markov Model (HMM)*: An algorithm that is used to decipher spoken words.

*Hyperparameters*: Features in a model that cannot be learned directly from the training process.

*Instance*: This is a row of data.

*Jupyter Notebook*: A web-based app that makes it easy to code in Python and R to create visualizations and import AI systems.

*K-Means Clustering*: An algorithm that is effective for grouping similar unlabeled data.

*K-Nearest Neighbor (k-NN)*: A machine learning algorithm that classifies data based on similarities.

*Lemmatization*: A process in NLP that removes affixes or prefixes so as to focus on finding similar root words.

*Lidar (Light Detection and Ranging):* A device—which is usually at the top of an autonomous car—that shoots laser beams to measure the surroundings.

*Linear Regression*: Shows the relationship between certain variables, which can help with predictions for machine learning systems.

*Machine Learning*: Where a computer can learn and improve by processing data without having to be explicitly programmed. Machine learning is a subset of AI.

*Metadata*: This is data about data—that is, descriptions. For example, a music file can have metadata like the size, length, date of upload, comments, genre, artist, and so on.

*Naïve Bayes Classifier*: A method of machine learning that uses Bayes' theorem to make predictions, but the variables are independent from each other.

*Named Entity Recognition*: In the NLP process, this involves identifying words that represent locations, persons, and organizations.

*Natural Language Processing (NLP)*: A subset of AI that deals with how computers understand and manipulate language.

*Neural Network*: A sophisticated AI model that mimics the brain. A neural network has various layers that attempts to find unique patterns that involve multiple layers of analysis.

*Normal Distribution*: A plot of data that looks like a bell and the midpoint is the mean.

*NoSQL System*: A next-generation database. The information is based on a document model so as to allow for more flexibility with analysis as well as the handling of structured and unstructured data.

*Ordinal Data*: A mix of numerical and categorical data, such as an Amazon.com rating for a product.

*Overfitting*: Where a model is not accurate because the data is not reflective of what is being tested or there is a focus on the wrong features.

*Pearson Correlation*: Shows the strength of a correlation—from 1 to -1. The closer it is to 1, the more accurate the correlation.

*Phonemes*: The most basic units of sound in a language.

*Predictive Analytics*: Involves using data to make forecasts.

*Python*: A computer language that has become the standard in developing AI models.

*PyTorch*: A platform, developed by Facebook, that allows for the creation of sophisticated AI models.

*Recurrent Neural Network (RNN)*: A deep learning model that processes prior inputs across time. A common use case is when a person types in characters in a messaging app, as the AI will predict the next word.

*Reinforcement Learning*: An approach to creating an AI model where the system is rewarded for the right predictions and punished for the wrong ones.

*Relational Database*: A database, whose roots go back to the 1970s, that creates relationships among tables of data and has a scripting language, called SQL.

*Robotic Desktop Automation (RDA):* The RPA system works alongside an employee to handle jobs or tasks.

*Robotic Process Automation (RPA)*: A category of software that automates routine and mundane tasks within an organization. It is often an initial way to implement AI.

*Robot Operating System (ROS)*: An open source middleware system that manages critical parts of a robot.

*R-squared*: Provides a way to gauge the accuracy of a regression. An R-squared ranges from 0 to 1. And the closer a model is to 1, the higher the accuracy.

*Sentiment Analysis*: This is where you mine social media data and find the trends.

*Sensor*: The typical sensor is a camera or a Lidar, which uses a laser scanner to create 3D images.

*Sigmoid*: A common activation function for a deep learning model. It has a value that ranges from 0 to 1. What's more, the closer it is to 1, the higher the accuracy.

*Standard Deviation*: Measures the average distance from the mean, which gives a sense of the variation in the data.

*Stemming*: Describes the process of reducing a word to its root (or lemma), such as by removing affixes and suffixes.

*Strong AI*: This is true AI, in which a machine is able to engage in human-like abilities like open-ended discussions.

*Structured Data*: Data that is usually stored in a relational database or spreadsheet, as the information is in a preformatted structure (like Social Security numbers, addresses, and point of sale information).

*Tagging Parts of Speech (POS)*: In the NLP process, this involves going through text and designating each word to its proper grammatical form, say nouns, verbs, adverbs, etc.

*TensorFlow*: An open source platform, backed by Google, that allows for the creation of sophisticated AI models.

*Test Data*: Data that a model's accuracy is evaluated.

*Three Laws of Robotics*: Based on the science fiction writings of Isaac Asimov, these laws provide the basic framework for how robots should interact with society.

*Tokenization*: In the NLP process where text is parsed and segmented into various parts.

*Topic Modelling*: In the NLP process, this involves looking for hidden patterns and clusters in the text.

*Training Data*: Data that is used to create an AI algorithm.

*True Positive*: When a model makes a correct prediction.

*Turing Test*: Created by Alan Turing, this is a way to determine if a system has achieved true AI. The test involves a person who asks questions to two participants—one human, the other a computer. If it is not clear who is the human, then the Turing Test has been passed.

*Unattended Robotic Process Automation (RPA)*: The RPA system is completely autonomous as the bot runs in the background.

*Unstructured Data*: Data that does not have predefined formatting, such as images, videos, and audio files.

*Supervised Learning*: An AI model that uses labeled data. This is the most common approach.

*Unsupervised Learning*: Involves an AI model that uses unlabeled data. Generally, this means there will need to be deep learning systems to detect patterns.

*Vanishing Gradient Problem*: Explains how the accuracy decays as a deep learning model gets larger.

*Virtual Assistant*: An AI device that helps a person with his or her daily activities.

*Weak AI*: This is where AI is used for a particular use case, such as with Apple's Siri.

# Index